

# ***In silico* prediction of mutant HIV-1 proteases cleaving a target sequence**

Jan H. Jensen,<sup>1</sup> Martin Willemoës,<sup>2</sup> Jakob R. Winther,<sup>2</sup> Luca De Vico<sup>1,\*</sup>

**1** Department of Chemistry, University of Copenhagen, Universitetsparken 5, DK-2100 Copenhagen, Denmark

**2** Department of Biology, University of Copenhagen, Ole Maaløes Vej 5, DK-2200 Copenhagen, Denmark

\* Corresponding Author, Email: luca@chem.ku.dk

## **Abstract**

HIV-1 protease represents an appealing system for directed enzyme re-design, since it has various different endogenous targets, a relatively simple structure and it is well studied. Recently Chaudhury and Gray (Structure (2009) 17: 1636 – 1648) published a computational algorithm to discern the specificity determining residues of HIV-1 protease. In this paper we present two computational tools aimed at re-designing HIV-1 protease, derived from the algorithm of Chaudhuri and Gray. First, we present an energy-only based methodology to discriminate cleavable and non cleavable peptides for HIV-1 proteases, both wild type and mutant. Secondly, we show an algorithm we developed to predict mutant HIV-1 proteases capable of cleaving a new target substrate peptide, different from the natural targets of HIV-1 protease. The obtained *in silico* mutant enzymes were analyzed in terms of cleavability and specificity towards the target peptide using the energy-only methodology. We found two mutant proteases as best candidates for specificity and cleavability towards the target sequence.

## **List of Abbreviations**

**PR** HIV-1 protease.

**WT-PR** Wild type HIV-1 protease.

**mutant PR** Mutant HIV-1 protease.

**Pr3 set** Set of mutant proteases derived from Pr3, a mutant protease developed by Alvizo *et al.* These were heterodimer proteases.

**DR set** Set of mutant proteases derived as a subset of HIV-1 proteases that have been found to be drug resistant. These were homodimer proteases.

## Introduction

Proteases represent a class of enzymes ubiquitous in all living organisms, with multiple applications in industry and biotechnology research [1–3]. There is thus interest in designing new proteases capable of cleaving specific peptide sequences [4]. HIV-1 protease (PR) represents an attractive starting structure for directed enzyme re-design, since it is known to cleave a variety of sequences. PR is the enzyme responsible for processing the gag – pol fusion polyproteins of the HIV virus [5]. PR is an aspartic protease [6–8] and is a homodimer where each chain is composed of 99 residues. Wild type PR (WT-PR) is very specific for the endogenous cleavage sequences of the polyprotein (endogenous substrate peptides, Table 1S), even if the source of this specificity is still not completely clear. A series of other non-endogenous peptides have also been found to be cleaved by PR. The latest hypothesis on the origin of this specificity, called dynamic substrate envelope [9,10], states that peptides fitting into the protease cavity through a certain number of hydrogen bonds will be bound and possibly cleaved nearly regardless of their amino acid composition. In fact, there is no clear trend in amino acid sequence (e.g. a negatively charged amino acid in position P1 or a hydrophobic one in position P2’). This suggests that with few mutations PR could be made to cleave other target peptide sequences in a specific manner.

Many computational studies on PR, both wild type (WT) and drug resistant mutant enzymes, are aimed at elucidating the affinity of the enzymes towards endogenous substrates and inhibitors to be used as drug candidates [11–14]. Recently Chaudhury and Gray [15] published a computational algorithm specifically tailored for PR and aimed at the identification of the specificity determining residues. The algorithm is based on PyRosetta [16], a python script-based interface to Rosetta [17]. Thanks to the algorithm the authors were able to predict accurate protease – substrate complex structures (within 1.1 Å rms of the corresponding crystal structure) and introduced an energetic discrimination of cleavable peptides. More recently Alvizo *et al.* [18] employed computational methods to re-engineer a mutant PR (Pr3) more specific for one of the endogenous peptide sequences over two others.

The first aim of this study is to develop an energy-only based methodology to discern cleavable and non cleavable peptides for PRs, WT and mutant. This methodology is based on the qualitative evaluation of PR:peptide complexes binding energies and is derived from the algorithm developed by Chaudhury and Gray. The second aim is to search and define an algorithm to predict mutant PRs capable of cleaving a specific target peptide sequence different from any endogenous substrate. We use our cleavability discerning methodology on the suggested mutant proteases, in order to define the best guess in terms of specificity towards the peptide sequence. In other words, the sought after mutant structure has to show better and worse binding towards the target and endogenous peptides, respectively, than WT-PR. To the best of our knowledge ours is the first study aimed at predicting a mutant PR capable of cleaving specifically a non endogenous peptide sequence.

The paper is organized as follows: first, we present our computed binding energies for known cleavable and non-cleavable peptides bound to WT-PR, selected peptides bound to a set of single, double and triple mutants (Pr3 set) derived from Pr3 as developed by Alvizo *et al.*, and a set of known mutant PRs and peptide derived from drug resistance (DR set) studies [19–21]. Secondly, we present two different versions of our algorithm to determine mutant PRs that will cleave the sequence HFLSF\*MAIP, where the \* symbol indicates the desired cleaving site. A discussion about the best strategy to suggest mutant enzymes follows. The conclusions summarize the main findings of the paper, followed by a detailed description of the employed computational methods.

## Results and Discussion

### Development of a Cleavability Test

In general, the activity of an enzyme towards two similar substrates is regulated by (i) the strength of the enzyme-substrate binding and (ii) the efficiency of the enzymatic reaction. The two processes are regulated by two constants, usually indicated as  $k_m$  and  $k_{cat}$ , respectively. The overall enzymatic efficiency is given by the ratio of these two constants. The dynamic substrate envelope hypothesis [10] suggests that if a peptide is bound to PR it will be cleaved. Thus, we decided to evaluate the binding energy of different peptides to PR, which can be correlated to  $k_m$ . We then compared the computed binding energies to PR of known cleavable and non cleavable peptides, to be correlated to corresponding ranges of binding energies. By so doing we disregarded  $k_{cat}$ , that is we did not consider possible effects from the enzymatic reaction.

The cleavability test was developed by considering binding energies of WT-PR with its endogenous and known cleavable substrates and known non-cleavable peptides. Afterwards we investigated the reliability of the test with mutant PRs (the Pr3 set) when binding PR endogenous substrates. Finally, we assessed the test on mutant PRs (the DR set) when binding mutant substrates.

The complete methodology for evaluating binding energies is described in the computational methods section. In brief, it is composed by a structure optimization algorithm, followed by an energetic re-evaluation of the obtained structures. In the following paragraphs we evaluate our methodology in terms of binding energies versus cleavability for: (1) WT-PR and its endogenous substrates and known cleavable and non-cleavable peptides, (2) the Pr3 set of mutant PRs and endogenous substrates and (3) the DR set with wild type and mutated endogenous peptides. Binding energies were computed also for WT-PR and all mutant PRs in complex with octa-alanine (poly-Ala) and octa-arginine (poly-Arg) peptides to test for aspecific binding.

(1) Table 1 reports the computed binding energies of the set of known cleavable endogenous peptides of WT-PR. The sequence of the tested endogenous cleavable peptides is reported in Table 1S. Alongside the endogenous peptides, a set of 59 known cleavable

peptides was also tested. The sequence of the 59 tested non-endogenous cleavable peptides was obtained as previously described [15, 22–26]. Table 5S reports the computed binding energies to WT-PR and Table 2S the sequences of these non-endogenous peptides. Table 6S reports the computed binding energies of a set of peptides supposedly non-cleavable by WT-PR. The sequence of the 43 tested non-cleavable peptides was obtained as previously described [15, 26, 27] and is reported in Table 3S.

We performed a Mann-Whitney’s U test [28] to compare the computed binding energies, and found a significant difference between the cleavable and non-cleavable sets ( $p \approx 10^{-7}$ ), as reported in Table 2. Thus, we deemed the binding energy criterion sufficient to achieve discrimination. We further analyzed the computed binding energies through an ROC plot [29] relative to different cutoff values, so as to differentiate between cleavable and non-cleavable peptides. The plot is reported in Figure 1, and the relative data in Table 11S. The computed area under ROC [30] is 0.79 and 0.80 for FMO and RosettaDock energies, respectively, being the values of 0.50 and 1.00 typical correspondingly of a useless and a perfect test. Through the ROC plot, we found the best cutoff values discerning cleavable and non-cleavable peptides as those closest to (0, 1), which represents the theoretical perfect test. We found that cutoff values of -25 kcal/mol and -3 kT are best at discerning FMO and RosettaDock computed binding energies, respectively. Both FMO and RosettaDock perform well in computing binding energies capable of discerning cleavable and non-cleavable peptides. However, Figure 1S shows that there is no apparent correlation between FMO and RosettaDock computed binding energies. Thus, we repeated the Mann-Whitney’s U test and ROC analysis excluding the set of non-endogenous known cleavable peptides binding energies. The rationale behind this analysis is that we expect WT-PR to bind the endogenous peptides with higher affinity, as opposed to the broader range of the complete cleavable set, characterized only by cleavability and not specificity. Consequently, we assume that the endogenous peptides set have better binding energies, than the complete set of cleavable peptides. The Mann-Whitney’s U test (Table 2) shows that the RosettaDock based binding energies are in this case two orders of magnitude worse than FMO at discerning cleavable and non-cleavable peptides. The relative ROC plot (Figure 2S) shows as well that the FMO data performs better than RosettaDock, in terms of more strict best cutoff value and larger area under the ROC. Thus, we concluded that FMO computed binding energies are better than RosettaDock ones since are capable of discerning expected effects, such as the usage of a better performing subset of peptides. In the rest of this paper we will discuss only binding energies computed through FMO energy re-evaluation.

From Table 1 it is expected that WT-PR exhibits qualitatively different binding to the poly-protein substrates, given their computed binding energies ranging from -41 for the binding of p6pol-PR to -72 kcal/mol for p2-NC, with an average value of -60 kcal/mol. However, available experimental  $K_m$  values [22] do not show any trend similar to the computed data. Still, one has to remember that these computed binding energies should be considered only qualitatively and only compared to others obtained in the same man-

ner. See the Computational Methods section for further details. Furthermore, the span of both computed energies for which experimental data are available (20 kcal/mol) and the  $K_m$  values (2 orders of magnitude) is too small to allow a clear trend. The computed binding energies for the set of cleavable non-endogenous peptides span a wide range of values, from -2 to -86 kcal/mol, with average -40 kcal/mol. These peptides not being the natural target of WT-PR may account for this large span. The average computed binding energy for all cleavable peptides is -43 kcal/mol. The computed binding energies for the non-cleavable set of peptides (Table 6S) span an even wider range of values than those of the cleavable ones. Some PR – peptides complexes show positive energies. The majority (56%) of the computed binding energies are in the range -35 – 0 kcal/mol. However, a few peptides show a binding energy to WT-PR similar to those of the cleavable peptides.

(2) Recently Alvizo *et al.* [18] suggested through computational means a triple mutant (Pr3) with increased binding capability towards the endogenous Rtp51-Rtp66 cleavage sequence peptide compared to that towards other two cleavage sequences CA-p2 and p2-NC. The efficiency of Pr3 in cleaving preferentially Rtp51-Rtp66 was later experimentally verified. Pr3 was made by tethering a mutated chain of protease (A28S, D30F, G48R) to a wild type one. For comparison with our predicted mutant PRs, Table 3 reports our computed binding energies for the Pr3 three-fold mutant, as well for simpler one- and two-fold mutant PRs derived from Pr3 (Pr3 set), as compared to WT. Note, however, that experimental data are available only for the three-fold mutant PR. In our calculations, Pr3 set carried mutations only on chain A, while still being formed by two separate chains. We expected to find that Pr3 computed binding would be stronger towards Rtp51-Rtp66, while weaker towards CA-p2 and p2-NC, compared to WT-PR. The computed binding energies of the Pr3 set show that the mutant enzymes often have higher affinity for the desired Rtp51-Rtp66 peptide compared to CA-p2 and p2-NC. Most notably the double mutant A28S/G48R has a stronger computed binding energy towards the target peptide than WT-PR, while lower for the other two endogenous substrates. The binding energy test indicates that A28S/G48R (for which there is no experimental data available) would have been a more successful mutation than Pr3. Nevertheless, the possibility of using the binding energy test with mutant PRs was found viable.

(3) Finally we decided to apply the binding energy test to series of mutant PRs binding mutant endogenous substrates. Thus, we evaluated the binding energies of drug resistant HIV-1 proteases towards wild type and mutant substrate peptides. It has been found that mutations of the cleavage sites are correlated to mutations of the protease, often leading to drug resistance. We analyzed the K436R and A431V mutations of the NC-p1 Gag substrate peptide cleavage sequence in relation to a series of single mutations and one double mutation of HIV-1 protease (DR set). It has been reported [19] that a K436R mutation increases resistance to protease inhibitor drugs when combined with I50V, I84V and I84V/L90M PR mutations, while the A431V mutation results in a more efficient PR regardless of other mutations. We expected that the more efficient mutant PR – mutant peptide combinations were also characterized by stronger binding energies. Table 4

reports the results of our binding energy test for the DR set. Our methodology indicates cleavability for all combinations of mutant PRs and mutated NC-p1 substrate peptides. While there are some fluctuations in the binding energies, no clear pattern arises that can be related to the experimental findings. Possibly, the increased efficiency of drug resistant mutant proteases towards mutated peptides is related to  $k_{cat}$ . As previously stated, the effects of this constant are not considered by the present approach. Nevertheless, the binding energy test was found suitable also for combinations of mutant PRs with any peptide.

## Prediction and Analysis of Mutant PRs

The second aim of this study was to develop a computational methodology for the design of a mutant PR. The sought after enzyme had to be capable of cleaving a new target substrate different from the endogenous ones. The obtained mutant PR should also be specific for the target peptide sequence compared to the endogenous peptides. The chosen sequence for the target peptide was HLSF\*MAIP, where the \* symbol indicates the desired cleaving site. The sequence was extracted from that of  $\kappa$ -casein. Once candidate mutant PRs were obtained, we employed the binding energy test to assess the enzymes cleaving capabilities. The possibility of an increase in cleaving capability towards the target substrate was asserted by differences in binding energy between WT-PR and mutant PRs. We evaluated the binding energies of mutant PRs in complex with the TF-PR peptide, used as a starting template (see the Computational Methods section), and the CA-p2 and p2-NC peptides (for selected mutant PRs) in order to test the specificity of our mutant PRs.

The mutant-generating algorithm is described in details in the Computational Methods section. Two main strategies (Strategy1 and Strategy2) were employed for generating mutant PRs. In Strategy1, the side chains of only the 6 residues previously indicated as specificity determining [15] were allowed to change. The analysis of the binding energies of the mutant PRs generated by Strategy1 found the enzymes insufficient to perform the desired scope. This prompted us to further develop the algorithm. In Strategy2, the side chains of 26 residues were allowed to change. See the Computational Methods section for further details on the residues choice. The analysis of the binding energies of these mutant PRs found some of the predicted enzymes to be adequate to cleave the desired target sequence.

Tables 7S and 8S in the Supporting material reports the Strategy1 mutant PRs (**M1** – **M16**) and their computed binding energies towards the target peptide and TF-PR, CA-p2 and p2-NC endogenous peptides. Among these mutant PRs, **M5** shows the strongest binding energy towards the target peptide. However it has to be noted that the computed binding energy of **M5** towards the TF-PR peptide (used as a starting template for all mutant enzymes) is also stronger with respect to WT. Possibly **M5** is simply a better generic binder. To verify this hypothesis we tested **M5** as a binder also for other two endogenous peptide sequences, CA-p2 and p2-NC. Compared to WT-PR, **M5** has weaker



binding energy for the former peptide, but equal for the latter. In conclusion, **M5** is not predicted to be more specific for the target sequence than for the endogenous peptides. Moreover, **M5** was not directly predicted through Strategy1, but as a homodimeric derivative of **M2**, which shows only a small improvement in binding of the target peptide. All other mutant PRs suggested by Strategy1, **M1** – **M4** and **M6** – **M16**, were found having a weak binding energy towards the target peptide, with some of them showing prominently positive binding energies. It can be concluded that Strategy1 is unsatisfactory at predicting a mutant PR with an increased and specific affinity towards the target peptide. This is possibly due to the fact that allowing only six residues to change is too strict a condition to achieve a suitable mutant PR.

Thus, we decided to further improve the mutation algorithm by including more residues among those that can be changed. The six generations of mutant PRs computed through our Strategy2 mutant algorithm are presented in Table 5. We refer to them as generations since at each macro step of the algorithm the lowest in energy (as computed with the standard RosettaDock energy function) structure was used as starting point for the next step. The sixth generation (**M23**) did not produce any new change with respect to the fifth (**M22**), and the algorithm was consequently terminated. For each generation the structure with the lowest absolute energy was further optimized. After generation 1 two mutant structures were chosen (**M17** and **M18**) since they are very close in energy (as evaluated with the RosettaDock energy function, data not shown) but relatively different as mutation sites. In addition, an extra mutant PR (**M24**) was generated as homodimer of **M22**. The computed binding energies of the Strategy2 mutant PRs (**M17** – **M24**) are shown in Table 6. All Strategy2 mutant PRs show a binding energy towards the target sequence two to four fold stronger than WT-PR, with **M17** displaying the strongest binding energy. However, as for **M5**, binding energies towards the template peptide TF-PR as well as CA-p2 and p2-NC are also stronger than WT. Possibly **M17** is also a good but generic binder. Through the subsequent generations of mutant proteases, at last **M22** shows a binding energy towards the target peptide more than three fold stronger than WT, while the computed binding energy towards the natural endogenous substrates is weaker than WT. Similar results were obtained for its homodimer **M24**. **M22** and **M24** show binding energies below the cutoff value of -25 kcal/mol, and thus represent the best candidates to be further studied experimentally.

We compared the structures of WT-PR and **M24** as optimized while binding the target peptide. Figure 2 reports the superimposed backbones of the two enzymes after structure alignment. The two computed structures are quite coincident. Hence, it is expected that **M24** should retain the main structural features of the wild type enzyme. We also tried to analyze the choice of changed residues. Figure 3 shows that the residues that were changed from WT-PR to **M24** are disposed all around the bound peptide. Figures 3S – 14S given as supporting material compare each residue that differs between WT-PR and **M24**, while bound to the target peptide. Although it is evident that the A28S substitution on chain A introduces a hydrogen bond between the residue and the side chain of the serine in the

peptide (Figure 3S), the other substitutions are less easily rationalized. On going study aims at elucidating the role of the other residues substitutions.

It is interesting to note that Strategy2 mutated only 7 out of the 26 residues that were set as mutable in the method. It is also worth noting that of the 7 residues (A28, D30, K45, I50, P81, V82, I84) suggested by Strategy2 in the various mutant generations, A28, K45, P81 are not included in the set of major mutations site of HIV-1 protease responsible for drug resistance [31], that is: D30, V32, M46, I47, G48, I50, I54, Q58, T74, L76 V82, N83, I84, N88, L90. A28, K45, and P81 together with I50 are also not included in the specificity determining residues set [15]. However, A28 was located by Alvizo *et al.* for the Pr3 mutant [18]. We envision Strategy2 also as a tool to locate those residues most involved in binding a given substrate peptide.

From the analysis of the different PRs, mutant and wild type, and their binding energies, it is worth to note that WT-PR has a certain affinity with the octa-arginine peptide. Its computed binding energy is at the limit to consider the octa-arginine peptide as cleavable by WT-PR. Possibly this relatively strong binding is given by very few interactions. Accordingly, the single D30F change on chain A, that is changing one negatively charged residue into an aromatic hydrophobic one, is able to drop the computed binding energy to 0, as shown in Table 4. The currently going analysis of the residue by residue interactions for the modified side chains will give further information also on this aspect of the binding of PR.

Finally, it is interesting to note that the algorithm is not always preserving amino acid side chain changes through the generations. For example, I84V on chain A is introduced in **M18** and kept in **M19**, **M20** and **M21**, but later reverted. Possibly, an isoleucin in position 84 is energetically more favorable, given the other side chain changes.

## Conclusions

In the first part of this study we developed a methodology to test the cleavability of a peptide by HIV-1 protease (Tables 1 and 6S), solely based on the binding energy between the enzyme and the substrate. The methodology can also be applied to mutant PRs, Table 3. The technique is based on a PyRosetta algorithm generating, iteratively, optimized structures, coupled with an energy re-evaluation at a higher level of theory (FMO/PCM MP2/6-31G(d)).

In the second part of this study, the optimization algorithm was extended to permit the stochastic change of the side chain of selected residues, in order to better bind a given target peptide sequence. The selected target peptide was required to be different from the endogenous peptides. The desired outcome was a mutant PR with stronger and weaker predicted binding energy for the target and endogenous peptides, respectively, compared to WT-PR. The mutant PRs **M22** and **M24** generated through Strategy2 exhibit such desired characteristics (Table 6). We analyzed the backbone structure of WT-PR and



**M24** and found no major differences, thus indicating that **M24** should retain the general structure features of wild type HIV-1 protease. Strategy2 algorithm is able to predict mutations outside the usual set of residues involved in drug resistance, possibly giving an ulterior insight into the binding process of HIV-1 protease.

Ongoing experimental studies will show if and how well **M22** and/or **M24** bind and cleave the target sequence. Our current experimental and computational studies are also aimed at analyzing **M24** mutations, residue by residue and in combination, and their possible role in binding the target sequence. It is our hope that the experimental tests will provide enough information to be used to further improve the mutant generating algorithm. If the combination of computational algorithm and experimental verification is successful it will maybe permit the design of mutant PRs specific for any given substrate peptide.

## Computational Methods

In general, the activity of an enzyme towards two similar substrates is regulated by (i) how good the enzyme-substrate binding is and (ii) how efficient the enzymatic reaction is. Following the dynamic substrate envelope hypothesis [9,10], we assume a correlation between the binding of different peptides to PR and cleavability of the former. Thus we compute qualitative binding energies, on the premise that lower binding energy equals better cleavability.

### Binding Energies

#### PyRosetta Algorithm

The structure of wild type (WT) HIV-1 protease in complex with different octa-peptides was optimized using PyRosetta 1.1, [16] a python script-based interface to Rosetta, [17] and the algorithm depicted in Figure 4. The algorithm is based on the flexible peptide-docking algorithm used by Chaudhury and Gray [15] to identify in WT HIV-1 protease the active-site residues mostly involved in the discrimination of cleavable and non-cleavable peptides. Following their algorithm, the HIV-1 protease – peptide complexes are represented in atomic resolution, as opposed to a coarse-grain representation. With respect to the algorithm described in [15], our algorithm (Figure 4) has a larger number of cycles ( $8 \times 4 \times 6 = 192$  compared to  $8 \times 12 = 96$ ), and more 'small' and 'shear' moves for the perturbation of both the side chain and the backbone atoms. The side chain conformations are further optimized through a repacking algorithm [32] and using the extended Dunbrack library [33,34]. The moves are applied to all residues of the substrate peptide plus a selected number of residues of the protease, with the following criterion: all residues inside a 5 Å distance from any atom of the substrate peptide, plus all the residues reported as active by Chaudhury and Gray [15], plus their  $\pm 1$  neighbours, plus if one residue is included

on only one chain it is made to be included in both. After the moves, an energy minimization step is performed, based on the Davidon-Fletcher-Powell method [35, 36]. Each structure is then accepted or rejected based on a Monte Carlo (MC) criterion depending on the standard RosettaDock energy function [32, 33, 37–39]. Along the optimization a temperature gradient was applied, from an initial value of  $kT = 3.0$  to 1.0, unless differently stated. 500 decoy structures were generated using 5 parallel algorithm runs, each producing 100 structures.

The main difference with the algorithm of [15] is that after the algorithm produced 500 decoy structures, the lowest in energy is chosen and used as a starting structure for another cycle of optimization. This process is repeated  $K$  times, until convergence. It was found that, after at least 5 cycles, the computed RosettaDock energy did not change between subsequent cycles as soon as all 5 parallel runs of a single cycle produced structures with the same energy. Consequently, in order to render as automatic as possible the algorithm, the fact that  $K > 5$  and that each parallel run produced, as best structure, a decoy with the same energy was taken as a mark for convergence. It was found that, on average, a value of  $K = 20$  was sufficient. As an example, Figure 5 reports the energy of WT-PR bound to TF-PR along the optimization. The points at each step corresponds to the RosettaDock energy of the lowest in energy decoy out of the 500 computed at that particular step. Such structure would then be used as starting point for the next cycle. At the end of the  $K$  cycles the lowest in energy decoy is chosen as the PyRosetta optimized structure.

The same algorithm was also used for the optimization of mutant HIV-1 proteases (*vide infra*), the octa-peptides alone, and the protease alone as apo-protein.

The starting structures were prepared from that of HIV-1 protease in complex with an inhibitor (PDB accession code 1HXB [40]), considered as apo-protein. In order to place the substrate peptide, the structure of a D25N deactivated protease in complex with the natural substrate peptide p2-NC (PDB accession code 1KJ7 [9]) was aligned with respect to the backbone atoms of the protease ( $RMS = 0.436 \text{ \AA}$ ). The starting structure was then composed using the apo-protein from 1HXB and the substrate peptide from 1KJ7. All subsequent protease-peptide complexes were created starting from this structure and mutating the peptide accordingly. See Table 1S, Table 3S and Table 4S for a complete list of the considered substrate peptides. Hydrogen atoms were added through the program Pymol [41].

## Further Structures Optimization and Energetic Re-evaluation

The position of the hydrogen atoms of each PyRosetta generated structure was optimized using Open Babel [42] with the MMFF94 [43–47] force field. The energy of each structure was finally re-evaluated at the higher level of theory ‘FMO2-MP2/6-31G(d)/PCM[1]’. Single point energy evaluations were carried out using the fragment molecular orbital (FMO) approximation [48, 49], as implemented in GAMESS [50]. Each FMO calculation

was carried out at the MP2 level of theory [51] with the 6-31G(d) basis set [52, 53] and the Polarizable Continuum Model (PCM) approximation [54, 55]. Pairs of fragments separated by more than two van der Waals radii were calculated using a Coulomb expression for the interaction energy and ignoring correlation effects (RESDIM=2.0 RCORSD=2.0 in \$FMO). The input files for the FMO calculations were prepared using the program FRAGIT [56].

## Binding Energies Evaluation

The re-evaluated energy of every optimized structure was used to compute the binding energy of PR with different substrate peptides. The binding energy ( $E_{Bind}$ ) of HIV-1 protease (wild type or mutated) and a peptide was evaluated with equation (1), where  $E_{Complex}$  is the energy of the complex,  $E_{APO}$  the energy of the protease optimized as apo-protein,  $E_{Pep}$  the energy of the optimized peptide.

$$E_{Bind} = E_{Complex} - (E_{APO} + E_{Pep}) \quad (1)$$

These binding energies can not be directly compared to experimental values, for which a much more complex and accurate methodology is required [57]. These energies were used only to qualitatively compare different PR – peptide combinations.

## Mutation Algorithm

A similar procedure as that described in Figure 4 was used to produce mutant HIV-1 proteases, possibly capable of cleaving a given peptide different from the endogenous substrate peptides. The general idea was to 'expose' the protease to a different peptide and allow some residues to change in order to accommodate it better. A target octapeptide was chosen: HLSF\*MAIP, where the \* symbol indicates the desired cleaving site. The peptide sequence was extracted from that of  $\kappa$ -casein.

The assumption behind the algorithm is that lowering the energy of the PR – peptide complex by changing the side chains of selected residues would decrease also the binding energy, thus increasing the cleavability.

Two different methodologies were designed to predict mutant PRs, Strategy1 and Strategy2. The Strategy1 mutation algorithm is depicted in Figure 6. Each optimization step corresponds to the algorithm of Figure 4. In the mutation steps (also based on the previous algorithm), the Dunbrack library of rotamers includes all rotamers of all amino acids, but only for a selected number of residues. The six specificity determining residues, as found by [15], are chosen to be altered. In other words, during the mutation step, whenever one of the alterable residues is being optimized, the random choice of a test rotamer is among all possible amino acids. In Scheme A alterations are allowed on all 6 residues on both chains, for a total of 12 alterable residues. Thus, side-chain perturbation and repacking rotamer choice is performed randomly selecting among 12 x

20 = 240 possible amino acids. In Scheme B only alterations on L76 and V82 of Chain A and D30, I47, G48, and I84 of Chain B are allowed, for a total of 6 alterable residues. In this case, side-chain perturbation and repacking rotamer choice is performed with a random selection among  $6 \times 20 = 120$  possible amino acids. Each mutation step took ca. 40 hours on 5 cpus to produce 500 decoys. The lowest energy decoy is then chosen as starting structure for the next step. The energy of the structure is evaluated with the standard RosettaDock energy function. The residue reference energy part of the energy function [32] takes into account also the differences between different amino acids. In other words, energy differences between two mutant structures originates solely from different side chain interactions rather than from a different number of atoms.

Both the mutation and the optimization steps were repeated  $K'$  and  $K$  times, respectively. The mutation cycles are considered converged once two following cycles do not introduce new mutations. Different values of  $K$  and  $K'$  were found necessary to reach convergence. After a series of mutation cycles ( $K' \geq 8$ ), a series of optimization cycles was performed ( $K \geq 8$ ), followed by another usually shorter mutation cycle ( $K' \leq 3$ ) and finally a short optimization cycle ( $K \leq 3$ ).

Among the naturally cleaved peptides, TF-PR (sequence SFNF\*PQIT) was chosen as a starting substrate peptide, since it is the most similar, in terms of conserved residues, to the target peptide (sequence HLSE\*MAIP). Consequently, the optimized structure of WT protease in complex with the TF-PR peptide was chosen as starting template. The substrate peptide sequence was altered one amino acid at the time, as reported in Table 10S. After each peptide alteration, a series of protease mutation and optimization cycles were performed. Once convergence was reached, a new peptide single amino acid change was introduced and the procedure repeated. Different mutant PRs were obtained from different runs by changing a few parameters, e.g. the initial temperature of the simulation. These parameters are specified in Table 7S. Some mutant PRs were also produced by 'exposing' the protease directly to the target peptide without prior intermediates (mutation Scheme F). This last process required a higher number of  $K'$  cycles ( $K' \geq 15$ ), but without having to cycle through one substrate peptide residue at the time.

All mutant PRs obtained through Strategy1 were heterodimers. By simply equalizing alterations on both chains a number of extra homodimer mutant PRs were also obtained. These structures were subsequently optimized as previously described.

In Strategy2 the number of residues allowed to change was increased in order to include all amino acids residing inside a 3 Å radius from the TF-PR peptide. In other words, we chose those residues with at least one atom that is distant at most 3 Å from any atom of the substrate peptide. The specificity determining residues were also included in the set of alterable amino acids, if not already present. The residues Asp25, Thr26 and Gly27 of both chains were excluded from the set, since they represent the catalytic triad [5]. The full set of 26 residues is reported in Table 9S. Thus, side-chain perturbation and repacking rotamer choice is performed randomly selecting among  $26 \times 20 = 520$  possible amino acids. The mutant PRs were generated using the target peptide directly

(Scheme F). Each mutation step took a bit more than 3 days on 5 cpus to produce 500 decoy structures. An initial temperature of 9 kT was used.  $K' = 6$  mutation cycles were performed. The lowest in energy decoy after each mutation step was subsequently optimized (two after the first step). The sixth mutation step did not introduce any new mutation in PR and the mutation cycle was stopped.

Also the mutant PRs obtained through Strategy2 were heterodimers. Only the homodimer of the final mutant PR was considered, see Table 5.

## Acknowledgments

Computational resources were provided by the Danish Center for Scientific Computing (DCSC). LDV acknowledges S. Chaudhury and J. J. Gray for fruitful discussions about PyRosetta, for providing the set of non cleavable peptides and a copy of their algorithm script. C. Steinmann is acknowledged for help with the program FRAGIT and the FMO based calculations.

## Supporting material

Supporting material available: Tables 1S – 11S, Figures 1S – 14S, a movie showing the three dimensional structure of WT-PR bound to the target peptide, with highlighted the residues that are changed in **M24**.

## References

1. Rao MB, Tanksale AM, Ghatge MS, Deshpande VV (1998) Molecular and Biotechnological Aspects of Microbial Proteases. *Microbiology and Molecular Biology Reviews* 62: 597-635.
2. Gupta R, Beg Q, Lorenz P (2002) Bacterial alkaline proteases: molecular approaches and industrial applications. *Applied Microbiology and Biotechnology* 59: 15-32.
3. Li Q, Yi L, Marek P, Iverson BL (2013) Commercial proteases: Present and future. *FEBS Letters* 587: 1155 - 1163.
4. van Beilen JB, Li Z (2002) Enzyme technology: an overview. *Current Opinion in Biotechnology* 13: 338 - 344.
5. Brik A, Wong CH (2003) HIV-1 protease: mechanism and drug discovery. *Org Biomol Chem* 1: 5-14.

6. Navia MA, Fitzgerald PMD, McKeever BM, Leu CT, Heimbach JC, et al. (1989) Three-dimensional structure of aspartyl protease from human immunodeficiency virus HIV-1. *Nature* 337: 615-620.
7. Wlodawer A, Miller M, Jaskolski M, Sathyanarayana B, Baldwin E, et al. (1989) Conserved folding in retroviral proteases: crystal structure of a synthetic HIV-1 protease. *Science* 245: 616-621.
8. Davies DR (1990) The Structure and Function of the Aspartic Proteinases. *Annual Review of Biophysics and Biophysical Chemistry* 19: 189-215.
9. Prabu-Jeyabalan M, Nalivaika E, Schiffer CA (2002) Substrate Shape Determines Specificity of Recognition for HIV-1 Protease: Analysis of Crystal Structures of Six Substrate Complexes. *Structure* 10: 369 - 381.
10. Özen A, Haliloğlu T, Schiffer CA (2011) Dynamics of Preferential Substrate Recognition in HIV-1 Protease: Redefining the Substrate Envelope. *Journal of Molecular Biology* 410: 726 - 744.
11. Wang W, Kollman PA (2001) Computational study of protein specificity: The molecular basis of HIV-1 protease drug resistance. *Proceedings of the National Academy of Sciences* 98: 14937-14942.
12. Sherman W, Tidor B (2008) Novel Method for Probing the Specificity Binding Profile of Ligands: Applications to HIV Protease. *Chemical Biology & Drug Design* 71: 387-407.
13. Perez MAS, Fernandes PA, Ramos MJ (2010) Substrate Recognition in HIV-1 Protease: A Computational Study. *The Journal of Physical Chemistry B* 114: 2525-2532.
14. Lemmon G, Kaufmann K, Meiler J (2012) Prediction of HIV-1 Protease/Inhibitor Affinity using RosettaLigand. *Chemical Biology & Drug Design* 79: 888-896.
15. Chaudhury S, Gray JJ (2009) Identification of Structural Mechanisms of HIV-1 Protease Specificity Using Computational Peptide Docking: Implications for Drug Resistance. *Structure* 17: 1636 - 1648.
16. Chaudhury S, Lyskov S, Gray JJ (2010) PyRosetta: a script-based interface for implementing molecular modeling algorithms using Rosetta. *Bioinformatics* 26: 689-691.
17. Kaufmann KW, Lemmon GH, DeLuca SL, Sheehan JH, Meiler J (2010) Practically Useful: What the Rosetta Protein Modeling Suite Can Do for You. *Biochemistry* 49: 2987-2998.



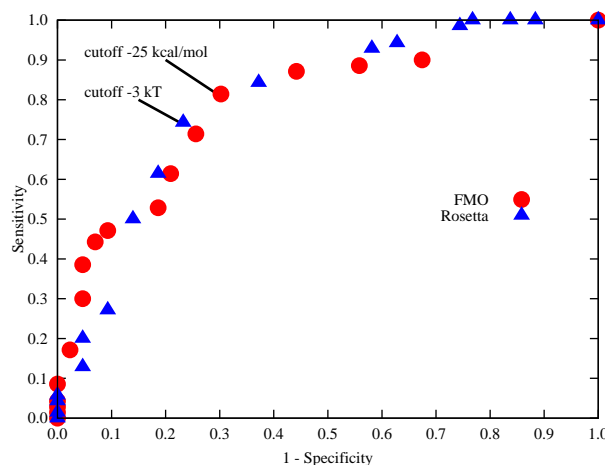
18. Alvizo O, Mittal S, Mayo SL, Schiffer CA (2012) Structural, kinetic, and thermodynamic studies of specificity designed HIV-1 protease. *Protein Science* 21: 1029–1041.
19. Kolli M, Stawiski E, Chappey C, Schiffer CA (2009) Human Immunodeficiency Virus Type 1 Protease-Related Cleavage Site Mutations Enhance Inhibitor Resistance. *Journal of Virology* 83: 11027-11042.
20. Özen A, Haliloğlu T, Schiffer CA (2012) HIV-1 Protease and Substrate Coevolution Validates the Substrate Envelope As the Substrate Recognition Pattern. *Journal of Chemical Theory and Computation* 8: 703-714.
21. Prabu-Jeyabalan M, Nalivaika EA, King NM, Schiffer CA (2004) Structural Basis for Coevolution of a Human Immunodeficiency Virus Type 1 Nucleocapsid-p1 Cleavage Site with a V82A Drug-Resistant Mutation in Viral Protease. *Journal of Virology* 78: 12446-12454.
22. Tözsér J, Bláha I, Copeland TD, Wondrak EM, Oroszlan S (1991) Comparison of the HIV-1 and HIV-2 proteinases using oligopeptide substrates representing cleavage sites in Gag and Gag-Pol polyproteins. *FEBS Letters* 281: 77 - 80.
23. Rivière Y, Blank V, Kourilsky P, Israël A (1991) Processing of the precursor of NF- $\kappa$ B by the HIV-1 protease during acute infection. *Nature* 350: 625 - 626.
24. Oswald M, von der Helm K (1991) Fibronectin is a non-viral substrate for the HIV proteinase. *FEBS Letters* 292: 298 - 300.
25. Tomaszek TA, Moore ML, Strickler JE, Sanchez RL, Dixon JS, et al. (1992) Proteolysis of an active site peptide of lactate dehydrogenase by human immunodeficiency virus type 1 protease. *Biochemistry* 31: 10153-10168.
26. Tomasselli AG, Sarcich JL, Barrett LJ, Reardon IM, Howe WJ, et al. (1993) Human immunodeficiency virus type-1 reverse transcriptase and ribonuclease h as substrates of the viral protease. *Protein Science* 2: 2167–2176.
27. Chou KC (1996) Prediction of Human Immunodeficiency Virus Protease Cleavage Sites in Proteins. *Analytical Biochemistry* 233: 1 - 14.
28. Mann HB, Whitney DR (1947) On a Test of Whether One of Two Random Variables is Stochastically Larger than the Other. *The Annals of Mathematical Statistics* 18: 50–60.
29. Metz CE (1978) Basic principles of {ROC} analysis. *Seminars in Nuclear Medicine* 8: 283 - 298.

30. Hanley JA, McNeil BJ (1982) The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 143: 29-36.
31. Johnson VA, Calvez V, Günthard HF, Paredes R, Pillay D, et al. (2013) Update of the Drug Resistance Mutations in HIV-1: March 2013. *Top Antivir Med* 21: 6-14.
32. Kuhlman B, Baker D (2000) Native protein sequences are close to optimal for their structures. *Proceedings of the National Academy of Sciences of the United States of America* 97: 10383-10388.
33. Dunbrack RL, Cohen FE (1997) Bayesian statistical analysis of protein side-chain rotamer preferences. *Protein Science* 6: 1661-1681.
34. Wang C, Schueler-Furman O, Baker D (2005) Improved side-chain modeling for protein-protein docking. *Protein Science* 14: 1328-1339.
35. Davidon WC (1991) Variable Metric Method for Minimization. *SIAM Journal on Optimization* 1: 1-17.
36. Fletcher R, Powell MJD (1963) A Rapidly Convergent Descent Method for Minimization. *The Computer Journal* 6: 163-168.
37. Gray JJ, Moughon S, Wang C, Schueler-Furman O, Kuhlman B, et al. (2003) Protein-Protein Docking with Simultaneous Optimization of Rigid-body Displacement and Side-chain Conformations. *Journal of Molecular Biology* 331: 281 - 299.
38. Lazaridis T, Karplus M (1999) Effective energy function for proteins in solution. *Proteins: Structure, Function, and Bioinformatics* 35: 133-152.
39. Kortemme T, Baker D (2002) A simple physical model for binding energy hot spots in protein-protein complexes. *Proceedings of the National Academy of Sciences of the United States of America* 99: 14116-14121.
40. Jaskolski M, Tomasselli AG, Sawyer TK, Staples DG, Heinrikson RL, et al. (1991) Structure at 2.5-Å resolution of chemically synthesized Human Immunodeficiency Virus Type 1 protease complexed with a hydroxyethylene-based inhibitor. *Biochemistry* 30: 1600-1609.
41. The PyMol Molecular Graphics System, Version 1.2r1 Schrödinger, LLC.
42. O’Boyle N, Banck M, James C, Morley C, Vandermeersch T, et al. (2011) Open Babel: An open chemical toolbox. *Journal of Cheminformatics* 3: 33.
43. Halgren TA (1996) Merck molecular force field. I. Basis, form, scope, parameterization, and performance of MMFF94. *Journal of Computational Chemistry* 17: 490-519.

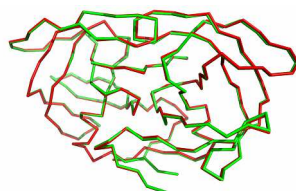
44. Halgren TA (1996) Merck molecular force field. II. MMFF94 van der Waals and electrostatic parameters for intermolecular interactions. *Journal of Computational Chemistry* 17: 520–552.
45. Halgren TA (1996) Merck molecular force field. III. Molecular geometries and vibrational frequencies for MMFF94. *Journal of Computational Chemistry* 17: 553–586.
46. Halgren TA, Nachbar RB (1996) Merck molecular force field. IV. Conformational energies and geometries for MMFF94. *Journal of Computational Chemistry* 17: 587–615.
47. Halgren TA (1996) Merck molecular force field. V. Extension of MMFF94 using experimental data, additional computational data, and empirical rules. *Journal of Computational Chemistry* 17: 616–641.
48. Nakano T, Kaminuma T, Sato T, Fukuzawa K, Akiyama Y, et al. (2002) Fragment molecular orbital method: use of approximate electrostatic potential. *Chemical Physics Letters* 351: 475 - 480.
49. Fedorov DG, Kitaura K (2007) Extending the Power of Quantum Chemistry to Large Systems with the Fragment Molecular Orbital Method. *The Journal of Physical Chemistry A* 111: 6904-6914.
50. Schmidt MW, Baldridge KK, Boatz JA, Elbert ST, Gordon MS, et al. (1993) General atomic and molecular electronic structure system. *Journal of Computational Chemistry* 14: 1347–1363.
51. Fedorov DG, Kitaura K (2004) Second order Møller-Plesset perturbation theory based upon the fragment molecular orbital method. *The Journal of Chemical Physics* 121: 2483-2490.
52. Hariharan PC, Pople JA (1973) The influence of polarization functions on molecular orbital hydrogenation energies. *Theoretical Chemistry Accounts: Theory, Computation, and Modeling (Theoretica Chimica Acta)* 28: 213-222.
53. Francl MM, Pietro WJ, Hehre WJ, Binkley JS, Gordon MS, et al. (1982) Self-consistent molecular orbital methods. XXIII. A polarization-type basis set for second-row elements. *The Journal of Chemical Physics* 77: 3654-3665.
54. Tomasi J, Mennucci B, Cammi R (2005) Quantum Mechanical Continuum Solvation Models. *Chemical Reviews* 105: 2999-3094.
55. Fedorov DG, Kitaura K, Li H, Jensen JH, Gordon MS (2006) The polarizable continuum model (PCM) interfaced with the fragment molecular orbital method (FMO). *Journal of Computational Chemistry* 27: 976–985.

56. Steinmann C, Ibsen MW, Hansen AS, Jensen JH (2012) FragIt: A Tool to Prepare Input Files for Fragment Based Quantum Chemical Calculations. PLoS ONE 7: e44480.
57. Genheden S, Kongsted J, Söderhjelm P, Ryde U (2010) Nonpolar Solvation Free Energies of Protein–Ligand Complexes. Journal of Chemical Theory and Computation 6: 3558-3568.

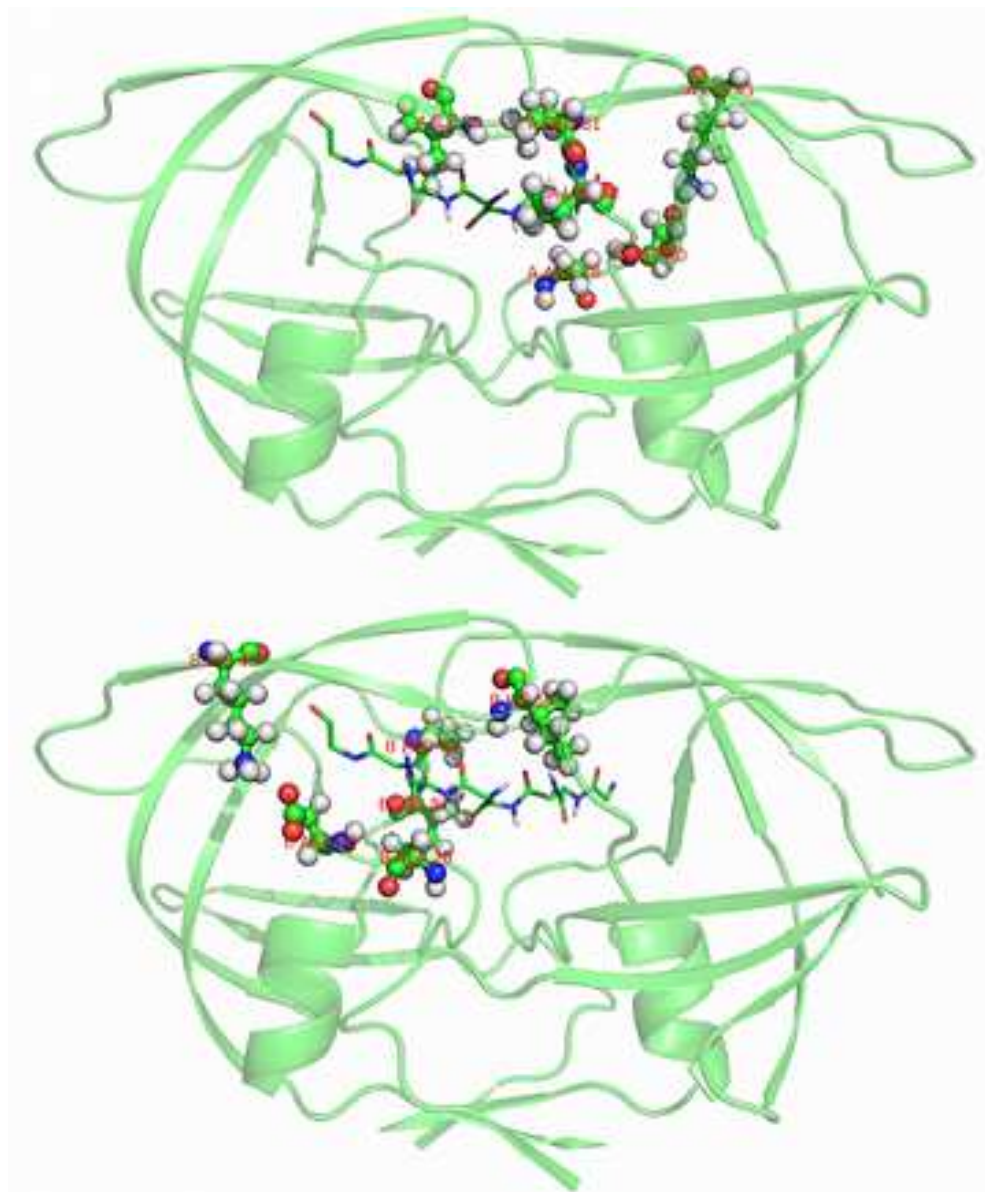
## Figure Legends



**Figure 1. ROC plot comparing different cutoff values for binding energies computed through FMO energy re-evaluation or RosettaDock energy function.** The values for each method closest to the theoretical optimum (0,1) are highlighted. The computed area under the ROC curve is 0.79 and 0.80 for FMO and Rosetta, respectively. The raw data is reported in Table 11S.

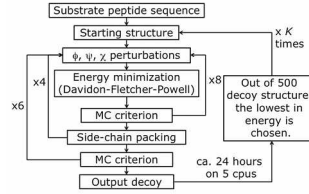


**Figure 2. Backbone difference between PyRosetta computed structures of WT-PR and M24.** The optimized structures of WT-PR and M24 binding the target peptide were aligned with respect to their  $\alpha$ -carbon atoms using PyMol. The backbone of M24 (red) is almost coincident with that of WT-PR (green) with a RMS of 0.227 Å.

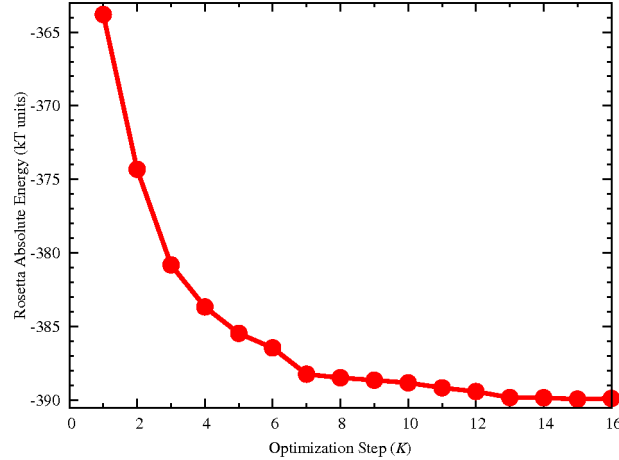


**Figure 3. Spatial disposition of the residues changed by Strategy2.** The six residues of chain A (top) and 6 residues of chain B (bottom) are highlighted in ball-and-sticks. The reported structure (as semi-transparent cartoon) is that of WT-PR optimized when binding the target peptide (only the backbone is shown in sticks). Figures 3S - 14S report the full residue by residue changes. A movie showing the three dimensional structure is included as Supporting Material.

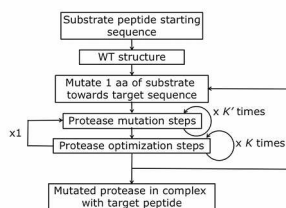




**Figure 4. PyRosetta based optimization algorithm.**  $\phi$ ,  $\psi$ ,  $\chi$  represent perturbations applied to both backbone and side chain dihedral angles. MC criterion stands for a Monte Carlo based check of decoy structures.



**Figure 5. Optimization algorithm convergence.** Example of energy convergence during the various macro cycles of the optimization algorithm for WT-PR in complex with TF-PR peptide. Each point along the graph corresponds to the energy (computed with the RosettaDock energy function) of the lowest in energy decoy out of 500 produced during each of the  $K$  steps.



**Figure 6. PyRosetta based mutation algorithm.** The optimization step is the algorithm presented in Figure 4. In the mutation step the side chain perturbation for the six specificity determining residues is among all possible rotamer of all 20 amino acids.

## Tables

**Table 1.** Computed binding energies of WT-PR and cleavable endogenous peptides.

Substrate Peptide	FMO (kcal/mol)	RosettaDock (kT)	exp $K_m$ (mM) [22]
MA-CA	-57	-9	0.15
CA-p2	-52	-6	0.01
p2-NC	-72	-4	0.05
NC-p1	-68	-3	
p1-p6	-47	1	
p6pol-PR	-41	-6	
TF-PR	-62	-5	<0.01
PR-RTp51	-64	-7	0.07
RTp51-RTp66	-68	-12	0.04
RTp66-INT	-62	-6	
RH-IN	-63	-10	0.006

**Table 2. Comparison of WT-PR computed binding energies.**

	RosettaDock Energy Function	FMO Energy Re-evaluation
Average endogenous <sup>a</sup>	-6 (kT)	-60 (kcal/mol)
(Standard deviation)	(3) (kT)	(10) (kcal/mol)
Average all cleavable <sup>b</sup>	-5 (kT)	-43 (kcal/mol)
(Standard deviation)	(3) (kT)	(22) (kcal/mol)
Average non cleavable <sup>c</sup>	-1 (kT)	-15 (kcal/mol)
(Standard deviation)	(4) (kT)	(28) (kcal/mol)
U test probability (all cleavable VS non-cleavable)	$1.46 \cdot 10^{-7}$	$2.21 \cdot 10^{-7}$
U test probability (only endogenous VS non-cleavable)	$3.49 \cdot 10^{-4}$	$6.16 \cdot 10^{-6}$

a Table 1.

b Table 1 plus Table 5S.

c Table 6S.

Energies were computed with the standard RosettaDock energy function, as described in [15] and with the FMO re-evaluation. A Mann-Whitney's U test probability was evaluated by comparing the binding energies of the set of endogenous peptide against the non-cleavable and the entire set of cleavable peptides against the non-cleavable. The FMO based binding energies are more clear in discriminating cleavable and non cleavable peptides than the Rosetta based ones.

**Table 3. FMO computed binding energies of HIV-1 protease WT and Pr3 set of mutant PRs.**

PR	Peptides					
	RTp51-RTp66	poly-Ala	poly-Arg	TF-PR	CA-p2	p2-NC
WT-PR	-68	-15	-41	-62	-52	-72
<i>Single mutant</i>						
A28S	-65	-12	-35	-41	-13	-67
D30F	-48	-1	0	-7	-4	-43
G48R	-66	-44	-27	-55	-15	-43
<i>Double mutant</i>						
A28SD30F	-44	-16	30	-35	-22	-55
A28SG48R	-96	-16	-54	-35	-14	-60
D30FG48R	-76	3	18	-19	-1	-54
<i>Triple mutant</i>						
A28SD30FG48R	-42	-21	12	-32	-12	-63

Computed binding energies (kcal/mol) of WT and mutant HIV-1 proteases in complex with RTp51-RTp66, poly-alanine, poly-arginine, TF-PR, CA-p2 and p2-NC peptides.

**Table 4. FMO computed binding energies of HIV-1 protease WT and selected drug resistance mutant PRs (DR set).**

PR	Peptides				
	NC-p1 <sup>WT</sup>	NC-p1 <sup>K436R</sup>	NC-p1 <sup>A431V</sup>	poly-Ala	poly-Arg
WT-PR	-49	-70	-56	-15	-41
D30N	-29	-44	-36	-15	-23
I50L	-64	-49	-54	-16	-49
I50V	-54	-46	-52	-18	-34
V82A	-45	-54	-52	-16	-38
I84V	-46	-75	-35	-23	-38
I84V L90M	-55	-67	-55	-20	-46

Computed binding energies (kcal/mol) of WT-PR and selected drug resistance mutant proteases in complex with NC-p1 as wild type, K436R and A431V drug resistance associated mutant peptides, poly-alanine and poly-arginine peptides.

**Table 5. Strategy2 suggested mutant PRs.**

Mutant ID	Chain A	Chain B	Mutation Scheme	Notes
<b>M17</b>	A28S D30T	A28S D30T	F	After one mutation step
		K45M I50L V82F		
<b>M18</b>	A28S D30T I50L	A28S D30T	F	After one mutation step
	P81D V82R I84V	K45M I50L V82Y		
<b>M19</b>	A28S D30T I50L	A28S D30T K45A	F	After two mutation steps
	P81D V82R I84V	I50L V82Y I84L		
<b>M20</b>	A28S D30T I50L	A28S D30T K45D	F	After three mutation steps
	P81D V82R I84V	I50L V82Y I84L		
<b>M21</b>	A28S D30T I50L	A28S D30T K45D	F	After four mutation steps
	P81L V82Y I84V	I50L V82Y		
<b>M22</b>	A28S D30T I50L	A28S D30T K45A	F	After five mutation steps
	P81L V82Y	I50L V82Y		
<b>M23</b>	A28S D30T I50L	A28S D30T K45A	F	After six mutation steps
	P81L V82Y	I50L V82Y		
<b>M24</b>	A28S D30T K45A	A28S D30T K45A	–	Homodimer of <b>M22</b>
	I50L P81L V82Y	I50L P81L V82Y		

**M17** – **M23** represent the subsequent generations of mutant PRs suggested by Strategy2. All mutant enzymes were generated following Scheme F.



**Table 6.** FMO computed binding energies of HIV-1 protease WT and Strategy2 mutant PRs.

PR	Peptides					
	Target	poly-Ala	poly-Arg	TF-PR	CA-p2	p2-NC
WT-PR	-9	-15	-41	-62	-52	-72
<i>Gen 1</i>						
<b>M17</b>	-34	-13	-47	-68	-82	-74
<b>M18</b>	-24	-19	-45	-82	-62	-63
<i>Gen 2</i>						
<b>M19</b>	-17	2	1	-67	-46	-81
<i>Gen 3</i>						
<b>M20</b>	-23	-2	-19	-67	-33	-84
<i>Gen 4</i>						
<b>M21</b>	-20	2	7	-37	-32	-18
<i>Gen 5</i>						
<b>M22</b>	-29	-6	10	-42	-25	-30
<b>M24</b>	-29	-11	7	-44	-33	-33

Computed binding energies (kcal/mol) of WT-PR and Strategy2 mutant proteases in complex with Target, poly-alanine, poly-arginine, TF-PR, CA-p2 and p2-NC peptides.

## Supporting Material for

# *In silico* prediction of mutant HIV-1 proteases cleaving a target sequence

Jan H. Jensen, Martin Willemoës, Jakob R. Winther, Luca De Vico

The video animation of the optimized structure of WT-PR binding the target peptide with highlighted residues can be found at this link: <http://youtu.be/NEXKojTw2Bc> .

**Table 1S. Cleavable peptides.**

	P4	P3	P2	P1	*	P1'	P2'	P3'	P4'
MA-CA	S	Q	N	Y	*	P	I	V	Q
CA-p2	A	R	V	L	*	A	E	A	M
p2-NC	A	T	I	M	*	M	Q	R	G
NC-p1	R	Q	A	N	*	F	L	G	K
p1-p6	P	G	N	F	*	L	Q	S	R
p6pol-PR	S	F	N	F	*	P	Q	V	T
TF-PR	S	F	N	F	*	P	Q	I	T
PR-RTp51	T	L	N	F	*	P	I	S	P
RTp51-RTp66	A	E	T	F	*	Y	V	D	G
RTp66-INT	R	K	V	L	*	F	L	D	G
RH-IN	R	K	I	L	*	F	L	D	G

List of the cleavable endogenous peptides considered in this work

Table 2S. Extra cleavable peptides.

	P4	P3	P2	P1	*	P1'	P2'	P3'	P4'
K001	T	Q	I	M	*	F	E	T	F
K002	G	Q	V	N	*	Y	E	E	F
K003	P	F	I	F	*	E	E	E	P
K005	D	T	V	L	*	E	E	M	S
K007	A	E	E	L	*	A	E	I	F
K008	S	L	N	L	*	R	E	T	Q
K010	A	E	C	F	*	R	I	F	D
K011	D	Q	I	L	*	I	E	I	C
K012	D	D	L	F	*	F	E	A	D
K013	Y	E	E	F	*	V	Q	M	M
K014	P	I	V	G	*	A	E	T	F
K016	R	E	A	F	*	R	V	F	D
K018	A	Q	T	F	*	Y	V	N	L
K019	P	T	L	L	*	T	E	A	P
K020	S	F	I	G	*	M	E	F	K
K021	D	A	I	N	*	T	E	F	K
K022	Q	I	T	L	*	W	Q	R	P
K023	E	L	E	F	*	P	E	G	G
K029	K	E	L	Y	*	P	L	T	S
K031	S	R	S	L	*	Y	A	S	S
K032	A	E	A	M	*	S	Q	V	T
K034	G	S	H	L	*	V	E	A	L
K035	G	G	V	Y	*	A	T	R	S
K036	F	R	S	G	*	V	E	T	T
K037	V	E	V	A	*	E	E	E	E
K038	L	P	V	N	*	G	E	F	S
K039	E	T	T	A	*	L	V	C	D
K040	H	L	V	E	*	A	L	Y	L
K041	H	Y	G	F	*	P	T	Y	G
K042	D	S	A	D	*	A	E	E	D
K043	G	W	I	L	*	G	E	H	G
K045	Q	A	I	Y	*	L	A	L	Q
K046	E	K	V	Y	*	L	A	W	V
K047	V	E	I	C	*	T	E	M	E
K048	T	Q	D	F	*	W	E	V	Q
K049	L	W	M	G	*	Y	E	L	H
K050	G	D	A	Y	*	F	S	V	P
K051	E	L	E	L	*	A	E	N	R
K052	S	K	D	L	*	I	A	E	I
K053	L	E	V	N	*	I	V	T	D
K054	I	I	V	A	*	C	E	G	N
K056	G	G	N	Y	*	P	V	Q	H
K057	A	R	L	M	*	A	E	A	L
K058	P	F	A	A	*	A	Q	Q	R
K059	P	R	N	F	*	P	V	A	Q
K060	G	L	A	A	*	P	Q	F	S
K061	S	L	N	L	*	P	V	A	K
K063	R	Q	V	L	*	F	L	E	K
K064	Q	M	I	F	*	E	E	H	G
SUB3	Q	I	T	L	*	W	K	R	P
T035	V	E	I	C	*	T	E	M	E
T084	T	Q	D	F	*	W	E	V	Q
T112	G	D	A	Y	*	F	S	V	P
T228	L	W	M	G	*	Y	E	L	H
T300	E	L	E	L	*	A	E	N	R
T322	S	K	D	L	*	I	A	E	I
T480	Q	A	I	Y	*	L	A	L	Q
T491	L	E	V	N	*	I	V	T	D
T529	E	K	V	Y	*	L	A	W	V

List of the cleavable non-endogenous peptides considered in this work

**Table 3S. Other peptides**

	P4	P3	P2	P1	*	P1'	P2'	P3'	P4'
Target	H	L	S	F	*	M	A	I	P
NC-p1 <sup>A431V</sup>	R	Q	V	N	*	F	L	G	K
NC-p1 <sup>K436R</sup>	R	Q	A	N	*	F	L	G	R
poli-Ala	A	A	A	A	*	A	A	A	A
poli-Arg	R	R	R	R	*	R	R	R	R

List of other peptides considered in this work

Table 4S. Non-cleavable peptides

	P4	P3	P2	P1	*	P1'	P2'	P3'	P4'
NBP1	V	N	C	A	*	K	K	I	V
NBP2	W	R	N	R	*	C	K	G	T
NBP3	M	M	K	S	*	R	N	L	T
NBP4	L	A	A	A	*	M	K	R	H
NBP5	T	T	Q	A	*	N	K	H	I
T015	G	M	D	G	*	P	K	V	K
T031	I	K	A	L	*	V	E	I	C
T033	A	L	V	E	*	I	C	T	E
T037	I	C	T	E	*	M	E	K	E
T039	T	E	M	E	*	K	E	G	K
T080	L	N	K	R	*	T	Q	D	F
T082	K	R	T	Q	*	D	F	W	E
T086	D	F	W	E	*	V	Q	L	G
T088	W	E	V	Q	*	L	G	I	P
T108	V	L	D	V	*	G	D	A	Y
T110	D	V	G	D	*	A	Y	F	S
T114	A	Y	F	S	*	V	P	L	D
T116	F	S	V	P	*	L	D	E	D
T224	E	P	P	F	*	L	W	M	G
T226	P	F	L	W	*	M	G	Y	E
T230	M	G	Y	E	*	L	H	P	D
T232	Y	E	L	H	*	P	D	K	W
T296	T	E	E	A	*	E	L	E	L
T298	E	A	E	L	*	E	L	A	E
T302	E	L	A	E	*	N	R	E	I
T304	A	E	N	R	*	E	I	L	K
T318	Y	Y	D	P	*	S	K	D	L
T320	D	P	S	K	*	D	L	I	A
T324	D	L	I	A	*	E	I	Q	K
T326	I	A	E	I	*	Q	K	Q	G
T441	Y	V	D	G	*	A	A	N	R
T476	K	T	E	L	*	Q	A	I	Y
T478	E	L	Q	A	*	I	Y	L	A
T482	I	Y	L	A	*	L	Q	D	S
T484	L	A	L	Q	*	D	S	G	L
T487	Q	D	S	G	*	L	E	V	N
T489	S	G	L	E	*	V	N	I	V
T493	V	N	I	V	*	T	D	S	Q
T495	I	V	T	D	*	S	Q	Y	A
T525	L	I	K	K	*	E	K	L	A
T527	K	K	E	K	*	V	Y	L	A
T531	V	Y	L	A	*	W	V	P	A
T533	L	A	W	V	*	P	A	H	K

List of non-cleavable peptides considered in this work

**Table 5S. Computed binding energies of WT-PR and non-endogenous cleavable peptides.**

Substrate Peptide	FMO (kcal/mol)	RosettaDock (kT)	Substrate Peptide	FMO (kcal/mol)	RosettaDock (kT)
K001	-22	-2	K043	-67	-5
K002	-33	-4	K045	-55	-5
K003	-16	-8	K046	-39	-4
K005	-24	-4	K047	-79	-5
K007	-27	-1	K048	-86	-9
K008	-40	-2	K049	-40	-5
K010	-55	-4	K050	-36	-7
K011	-42	-3	K051	-31	-2
K012	-22	-7	K052	-4	1
K013	-36	-2	K053	-51	-4
K014	-29	-3	K054	-56	-1
K016	-73	-3	K056	-70	-5
K018	-72	-4	K057	-7	-6
K019	-61	-9	K058	-35	-2
K020	-8	-3	K059	-58	-6
K021	-64	-8	K060	-36	-5
K022	-45	-5	K061	-34	-5
K023	-63	-11	K063	-61	-3
K029	-48	-4	K064	-2	-7
K031	-32	-5	SUB3	-30	-6
K032	-51	-6	T035	-67	-4
K034	-8	-8	T084	-81	-7
K035	-30	-8	T112	-69	-11
K036	-3	-3	T228	-29	-5
K037	-2	-4	T300	-29	-1
K038	-39	-5	T322	-21	1
K039	-10	0	T480	-62	-5
K040	-30	-6	T491	-56	-1
K041	-53	-6	T529	-27	1
K042	-29	-6			



**Table 6S.** Computed binding energies of WT-PR and non-cleavable peptides.

Substrate Peptide	FMO (kcal/mol)	RosettaDock (kT)	Substrate Peptide	FMO (kcal/mol)	RosettaDock (kT)
NBP1	-18	3	T296	43	2
NBP2	-21	3	T298	-54	-2
NBP3	-63	0	T302	-44	4
NBP4	-18	4	T304	1	-2
NBP5	-68	7	T318	23	0
T015	2	-3	T320	23	-2
T031	-29	-6	T324	-18	0
T033	-10	-2	T326	-10	5
T037	-28	-5	T441	-30	-2
T039	-36	-3	T476	-10	1
T080	-45	4	T478	-33	-2
T082	9	-2	T482	-45	-1
T086	-23	0	T484	-24	-2
T088	-42	3	T487	31	-3
T108	-12	-4	T489	-24	-6
T110	-21	-3	T493	-19	-1
T114	-16	-8	T495	-15	1
T116	-10	-2	T525	61	7
T224	48	-1	T527	-9	1
T226	-49	-9	T531	-5	-5
T230	19	-2	T533	-23	-7
T232	-13	-1			

**Table 7S. Strategy1 suggested mutant PRs.**

Mutant ID	Chain A	Chain B	Mutation Scheme	Notes
<b>M1</b>	V82R I84V	D30Y V82I	A	
<b>M2</b>	V82Y	D30V	B	
<b>M3</b>	D30T	D30V V82I	B F	
<b>M4</b>	D30Y V82R	D30Y V82R	–	Homodimer of <b>M1</b>
<b>M5</b>	D30V V82Y	D30V V82Y	–	Homodimer of <b>M2</b>
<b>M6</b>	D30V	D30V	–	Homodimer of <b>M3</b>
<b>M7</b>	D30T I84V	D30V V82F	A F	Initial temperature = 9 kT
<b>M8</b>		D30V	B F	Initial temperature = 9 kT
<b>M9</b>	D30T I47L L76F V82R I84T	D30E V82Y	A	Initial temperature = 6 kT
<b>M10</b>	V82Y	D30T I84L	B	Initial temperature = 6 kT
<b>M11</b>	D30V V82Y I84V	D30H I47L L76F V82Y	A	Initial temperature = 12 kT
<b>M12</b>	V82Y	D30T	B	Initial temperature = 12 kT
<b>M13</b>	D30E L76F V82R	D30E L76F V82R	–	Homodimer of <b>M9</b>
<b>M14</b>	D30T V82Y I84L	D30T V82Y I84L	–	Homodimer of <b>M10</b>
<b>M15</b>	D30H I47L V82Y	D30H I47L V82Y	–	Homodimer of <b>M11</b>
<b>M16</b>	D30T V82Y	D30T V82Y	–	Homodimer of <b>M12</b>

Different mutant PRs were obtained by small modifications of the mutation algorithm. Inside Strategy1 two different schemes were used when choosing which residues could mutate. In Scheme A all six specificity determining residues were allowed to mutate on both chains. In Scheme B only residues 76 and 82 were set as mutable on Chain A and 30, 47, 48, and 84 on Chain B. In addition, a straight forward variant of the algorithm was tested, as opposed to the step-wise one presented in Table 10S. In this variant (Scheme F) the protease was directly 'exposed' to the final target peptide sequence. A  $K$  value in the order of 20 was necessary. Other parameters that differ from those specified in the Computational Methods section are also highlighted.

**Table 8S. FMO computed binding energies of HIV-1 protease WT and Strategy1 mutant PRs.**

PR	Peptides Target	poly-Ala	poly-Arg	TF-PR	CA-p2	p2-NC
WT-PR	-9	-15	-41	-62	-52	-72
M1	-7	3	28	-41		
M2	-13	-9	12	-41		
M3	-8	-4	13	-43		
M4	-18	-24	17	-52		
M5	-30	-14	4	-54	-27	-73
M6	-14	2	-1	-14		
M7	-7	-2	3	-49		
M8	-10	-9	-5	-55		
M9	20	17	14	-36		
M10	40	11	-3	-43		
M11	5	10	45	-44		
M12	2	2	16	-43		
M13	3	2	-24	-52		
M14	3	8	17	-60		
M15	-9	2	71	-51		
M16	-8	-6	26	-44		

Computed binding energies (kcal/mol) of WT-PR and Strategy1 mutant proteases in complex with target, poly-alanine, poly-arginine, TF-PR, CA-p2 and p2-NC peptides.

**Table 9S. Residues set as mutable in Strategy 2.**

Chain A	Chain B
Arg 8	Arg 8
Ala 28	Leu 23
Asp 29	Ala 28
Asp 30	Asp 29
Val 32	Asp 30
Gly 48	Lys 45
Gly 49	Ile 47
Ile 50	Gly 48
Leu 76	Gly 49
Thr 80	Ile 50
Pro 81	Pro 81
Val 82	Val 82
Ile 84	Ile 84

The residues were selected as those inside a 3 Å radius from the substrate peptide plus the specificity determining residues, if not included, minus the catalytic triad Asp25, Thr26 and Gly27 on both chains. The optimized structure of WT protease in complex with the TF-PR peptide was used as template.

**Table 10S. Substrate peptide mutation sequence**

	P4	P3	P2	P1	*	P1'	P2'	P3'	P4'
Start	Ser	Phe	Asn	Phe	*	Pro	Gln	Ile	Thr
	His	Phe	Asn	Phe	*	Pro	Gln	Ile	Thr
	His	Phe	Asn	Phe	*	Pro	Gln	Ile	Pro
	His	Leu	Asn	Phe	*	Pro	Gln	Ile	Pro
	His	Leu	Asn	Phe	*	Pro	Ala	Ile	Pro
	His	Leu	Ser	Phe	*	Pro	Ala	Ile	Pro
Target	His	Leu	Ser	Phe	*	Met	Ala	Ile	Pro

Step wise sequence of substrate peptides employed in the mutation algorithm. The starting sequence corresponds to the natural substrate TF-PR. This sequence is altered one amino acid at the time towards that of the desired target sequence. The P1 and P3' position were not changed during the sequence.

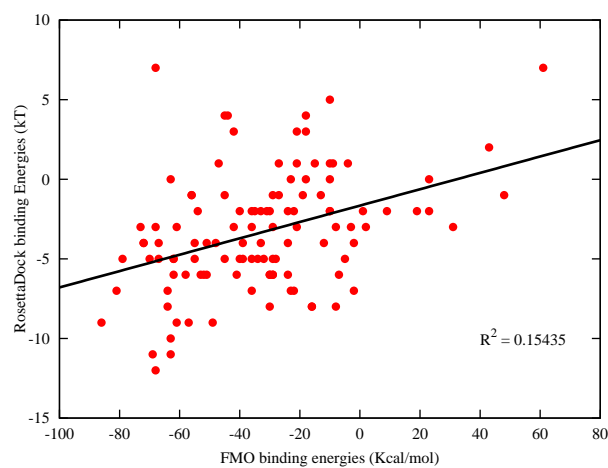
**Table 11S. ROC data.**

<i>Total</i>							
FMO				Rosetta			
cutoff	True positive	False positive	Distance to (0, 1)	cutoff	True positive	False positive	Distance to (0, 1)
-inf	1.00	1.00	1.00	-inf	1.00	1.00	1.00
-10	0.90	0.67	0.68	4	1.00	0.88	0.88
-15	0.89	0.56	0.57	3	1.00	0.84	0.84
-20	0.87	0.44	0.46	2	1.00	0.77	0.77
-25	0.81	0.30	0.35	1	0.99	0.74	0.74
-30	0.71	0.26	0.38	0	0.94	0.63	0.63
-35	0.61	0.21	0.44	-1	0.93	0.58	0.59
-40	0.53	0.19	0.51	-2	0.84	0.37	0.40
-45	0.47	0.09	0.54	-3	0.74	0.23	0.35
-50	0.44	0.07	0.56	-4	0.61	0.19	0.43
-55	0.39	0.05	0.62	-5	0.50	0.14	0.52
-60	0.30	0.05	0.70	-6	0.27	0.09	0.73
-65	0.17	0.02	0.83	-7	0.20	0.05	0.80
-70	0.09	0.00	0.91	-8	0.13	0.05	0.87
-75	0.04	0.00	0.96	-9	0.06	0.00	0.94
-80	0.03	0.00	0.97	-10	0.04	0.00	0.96
-85	0.01	0.00	0.99	-11	0.01	0.00	0.99
+inf	0.00	0.00	1.00	+inf	0.00	0.00	1.00

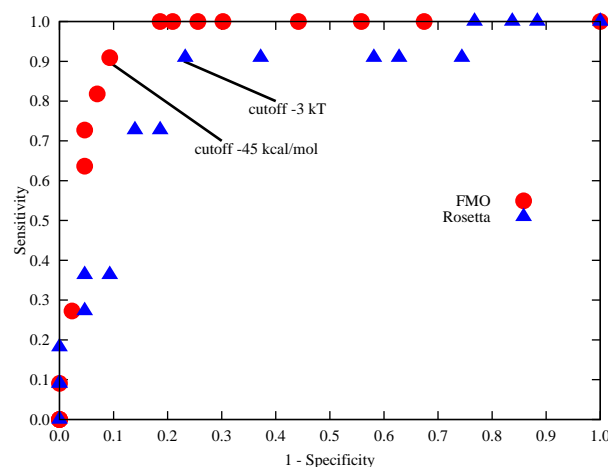
  

<i>Only endogenous</i>							
FMO				Rosetta			
cutoff	True positive	False positive	Distance to (0, 1)	cutoff	True positive	False positive	Distance to (0, 1)
-inf	1.00	1.00	1.00	-inf	1.00	1.00	1.00
-10	1.00	0.67	0.67	4	1.00	0.88	0.88
-15	1.00	0.56	0.56	3	1.00	0.84	0.84
-20	1.00	0.44	0.44	2	1.00	0.77	0.77
-25	1.00	0.30	0.30	1	0.91	0.74	0.75
-30	1.00	0.26	0.26	0	0.91	0.63	0.63
-35	1.00	0.21	0.21	-1	0.91	0.58	0.59
-40	1.00	0.19	0.19	-2	0.91	0.37	0.38
-45	0.91	0.09	0.13	-3	0.91	0.23	0.25
-50	0.82	0.07	0.19	-4	0.73	0.19	0.33
-55	0.73	0.05	0.28	-5	0.73	0.14	0.31
-60	0.64	0.05	0.37	-6	0.36	0.09	0.64
-65	0.27	0.02	0.73	-7	0.36	0.05	0.64
-70	0.09	0.00	0.91	-8	0.27	0.05	0.73
-75	0.00	0.00	1.00	-9	0.18	0.00	0.82
-80	0.00	0.00	1.00	-10	0.09	0.00	0.91
-85	0.00	0.00	1.00	-11	0.09	0.00	0.91
+inf	0.00	0.00	1.00	+inf	0.00	0.00	1.00

Comparison of ROC data for FMO energy re-evaluation and RosettaDock energy function generated binding energies, while considering different cutoff values. The upper part of the table reports the full comparison between known cleavable and non cleavable peptides. In the lower part, data for only the endogenous peptides was used for the cleavable part. True positive data is reported in the graphs as sensitivity, false positive as 1 - specificity. Theoretical values for  $\pm$  infinite cutoff have been added.



**Figure 1S. Correlation plot between FMO and RosettaDock computed binding energies.** The linear trend line shows no correlation between the data ( $R^2 = 0.15435$ ).



**Figure 2S. ROC plot comparing different cutoff values for binding energies computed through FMO energy re-evaluation or RosettaDock energy function.** The values for each method closest to the theoretical optimum (0,1) are highlighted. The comparison was done using the data of only the endogenous peptides for the cleavable peptides part. The computed area under the ROC curve is 0.96 and 0.84 for FMO and Rosetta, respectively. The raw data is reported in Table 11S.

Figures 3S - 14S compare the changes between WT-PR and **M24** , residue by residue. In each figure the enzyme is represented as semi-transparent ribbon, the peptide as sticks and the changing residue as ball-and-sticks. The peptide residues numbering is from 2 to 9. Each residue changed by Strategy2 is indicated by a label containing the chain, the residue name and its number.

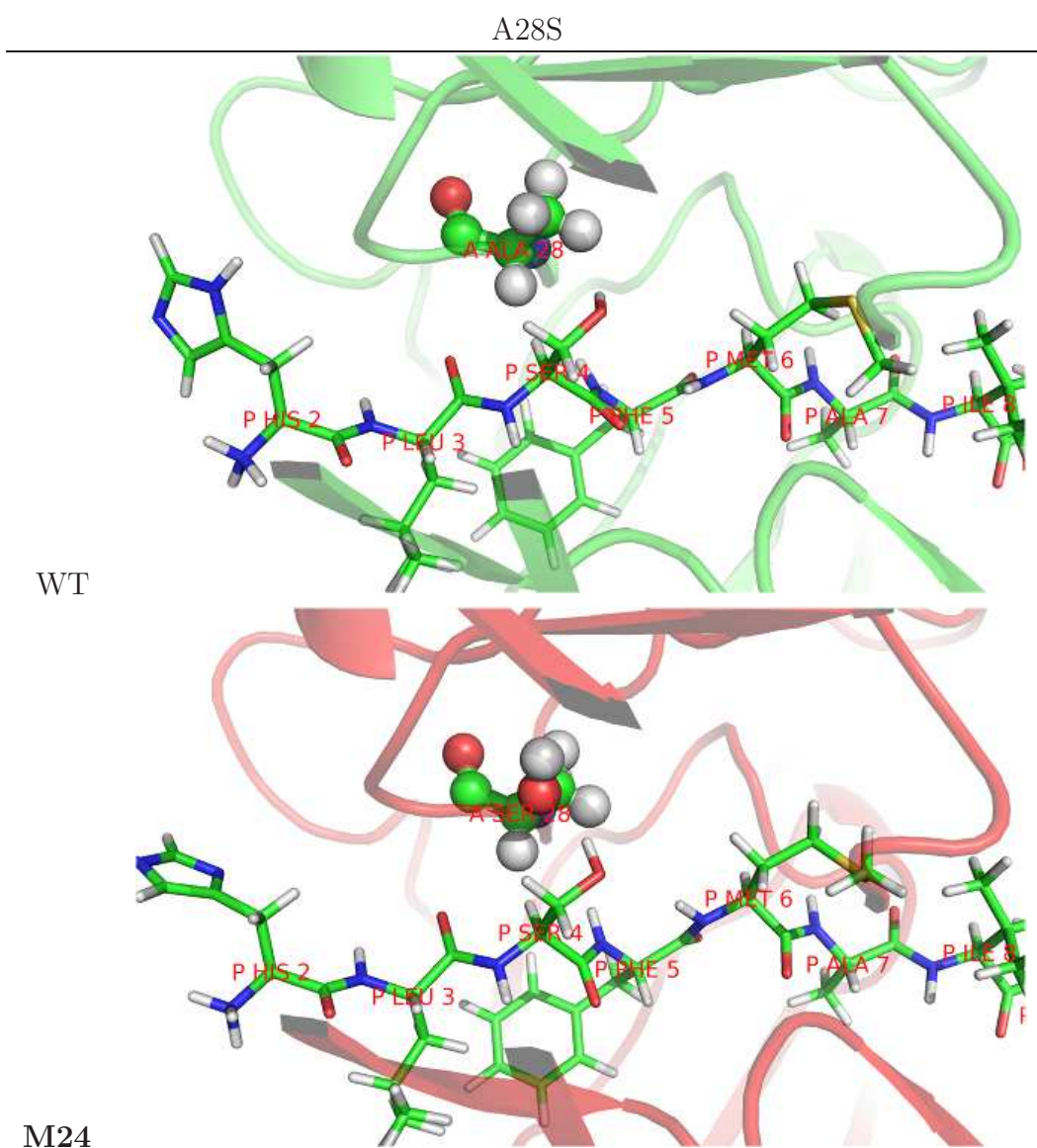


Figure 3S. Chain A, residue 28.



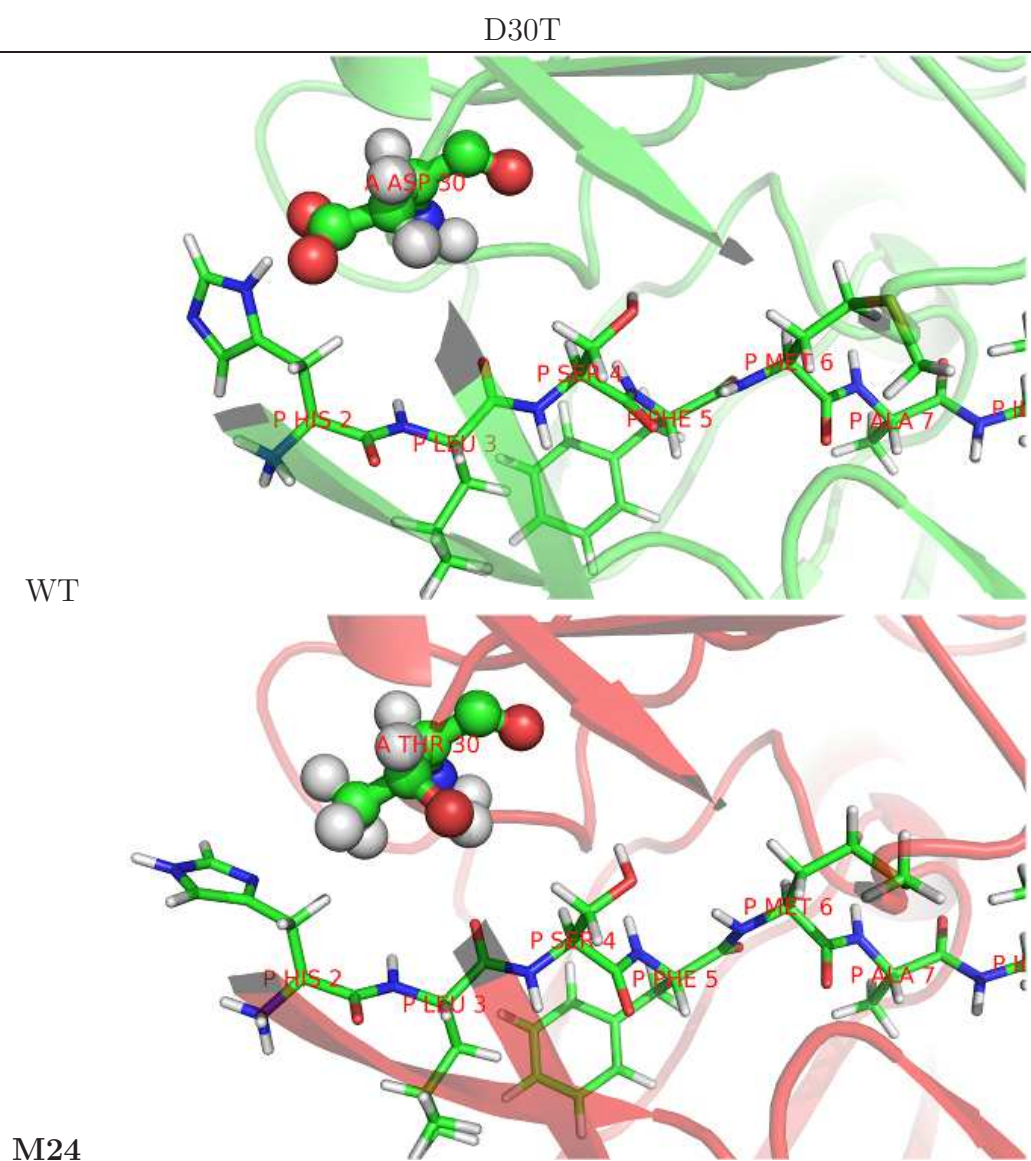


Figure 4S. Chain A, residue 30.

Figure 5S. Chain A, residue 45.

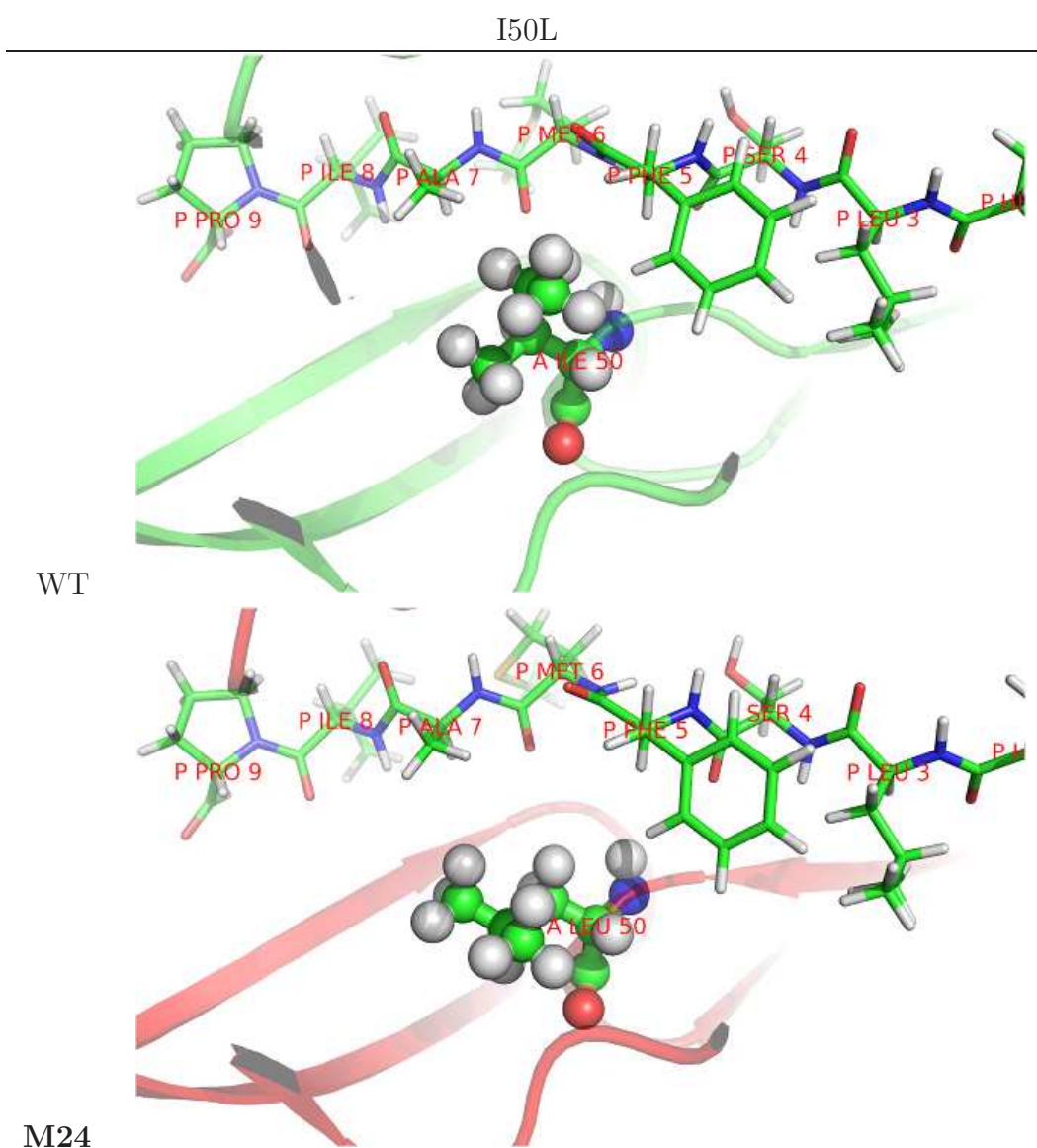


Figure 6S. Chain A, residue 50.

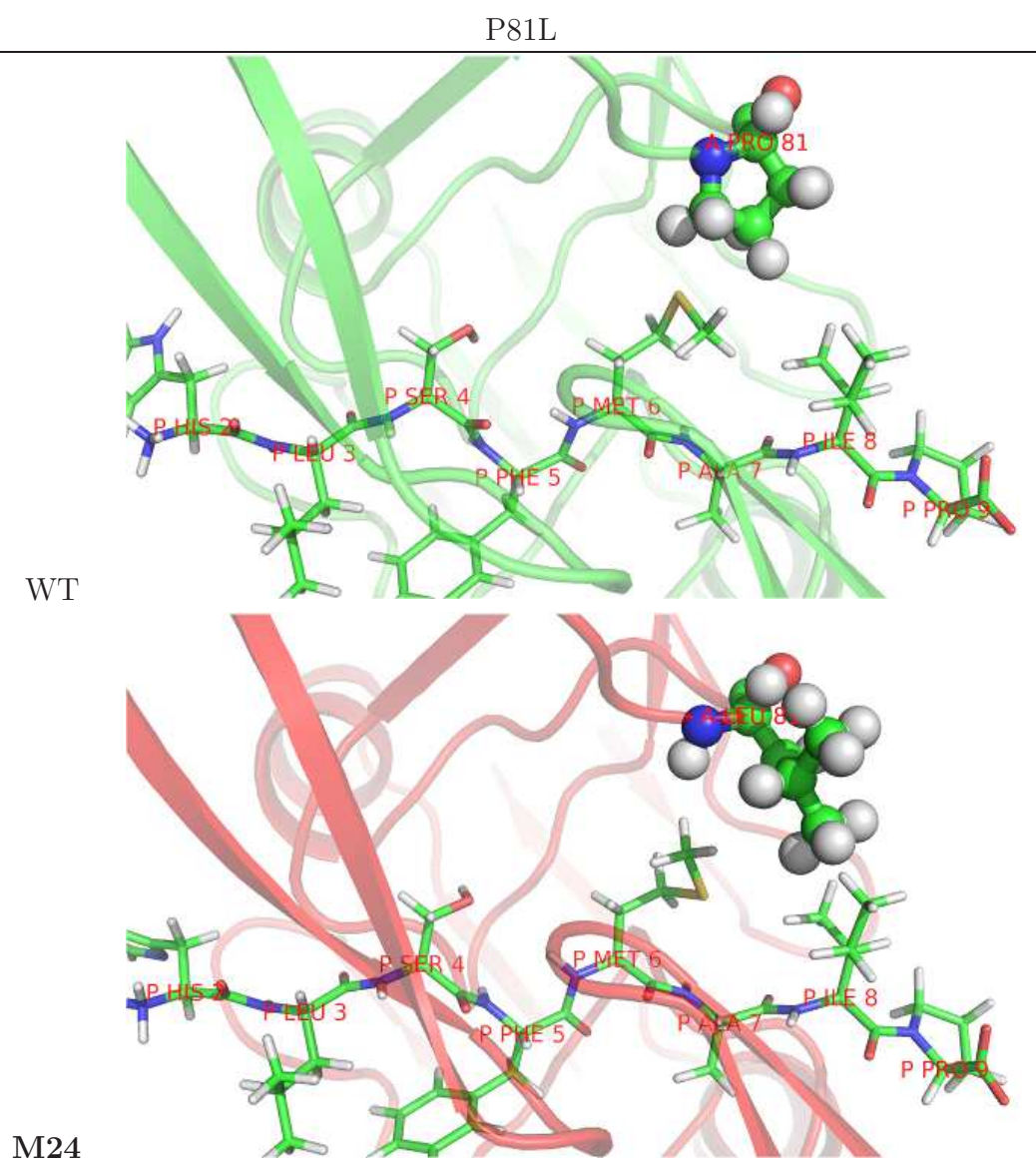


Figure 7S. Chain A, residue 81.



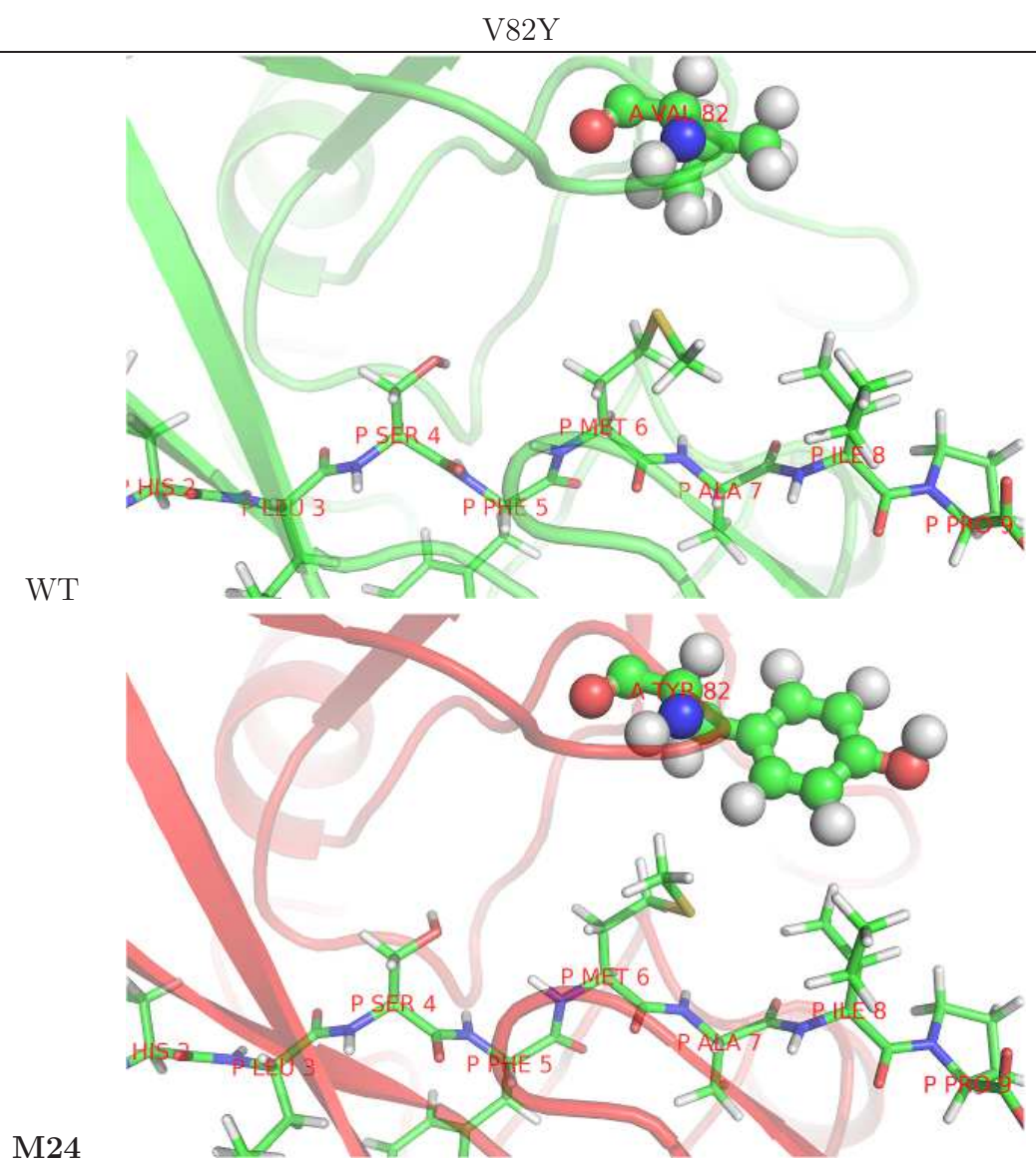


Figure 8S. Chain A, residue 82.

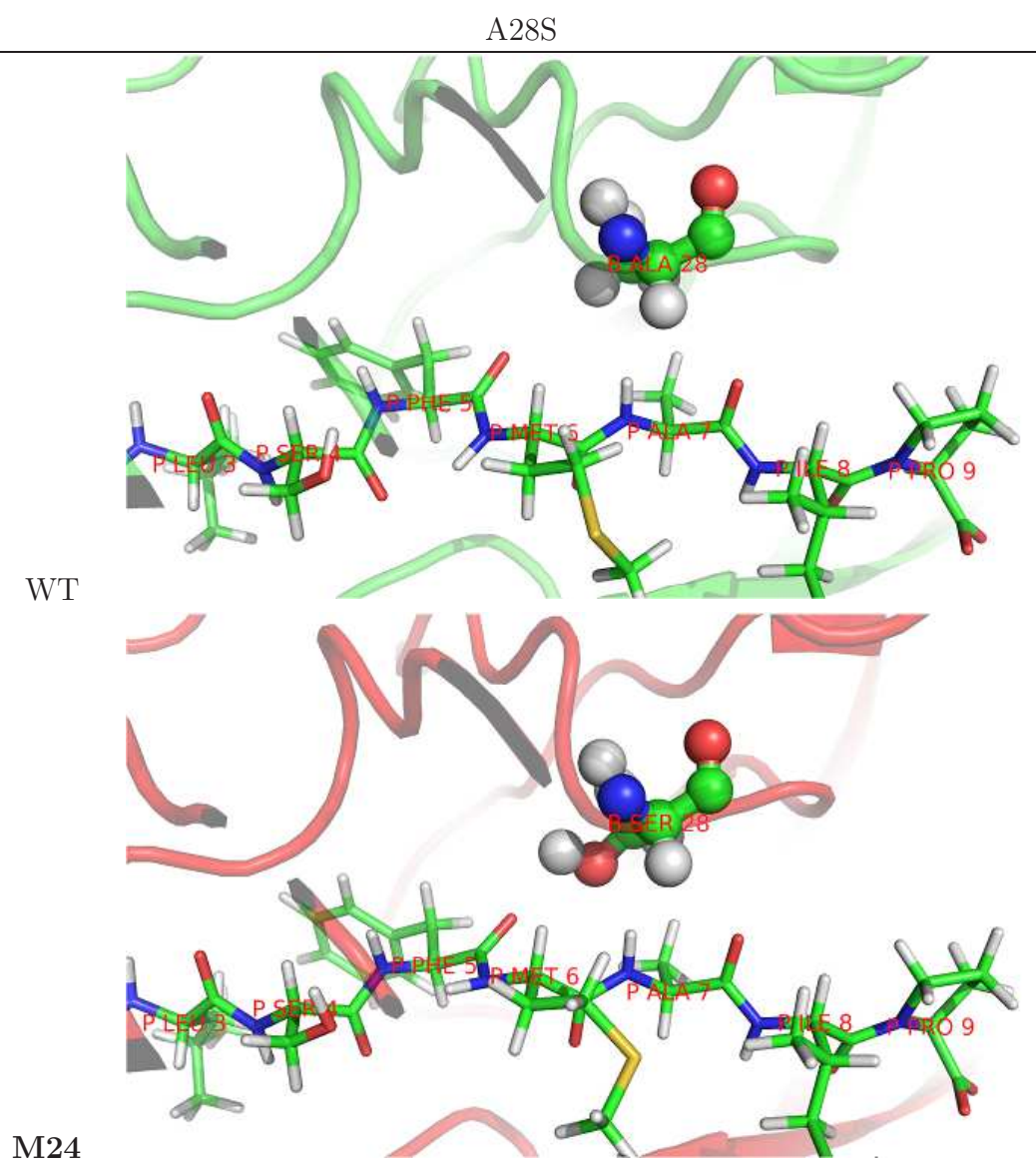


Figure 9S. Chain B, residue 28.

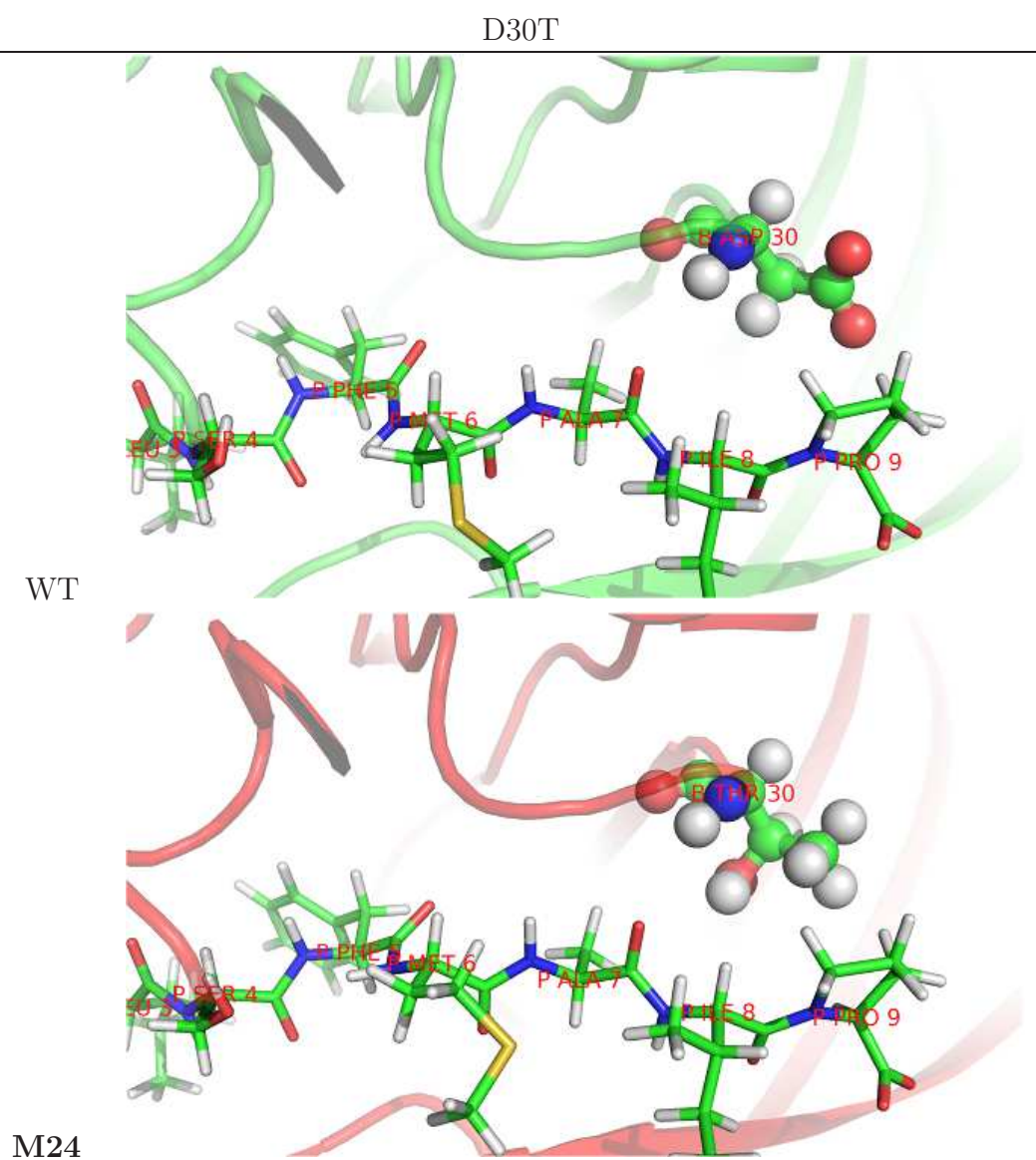


Figure 10S. Chain B, residue 30.

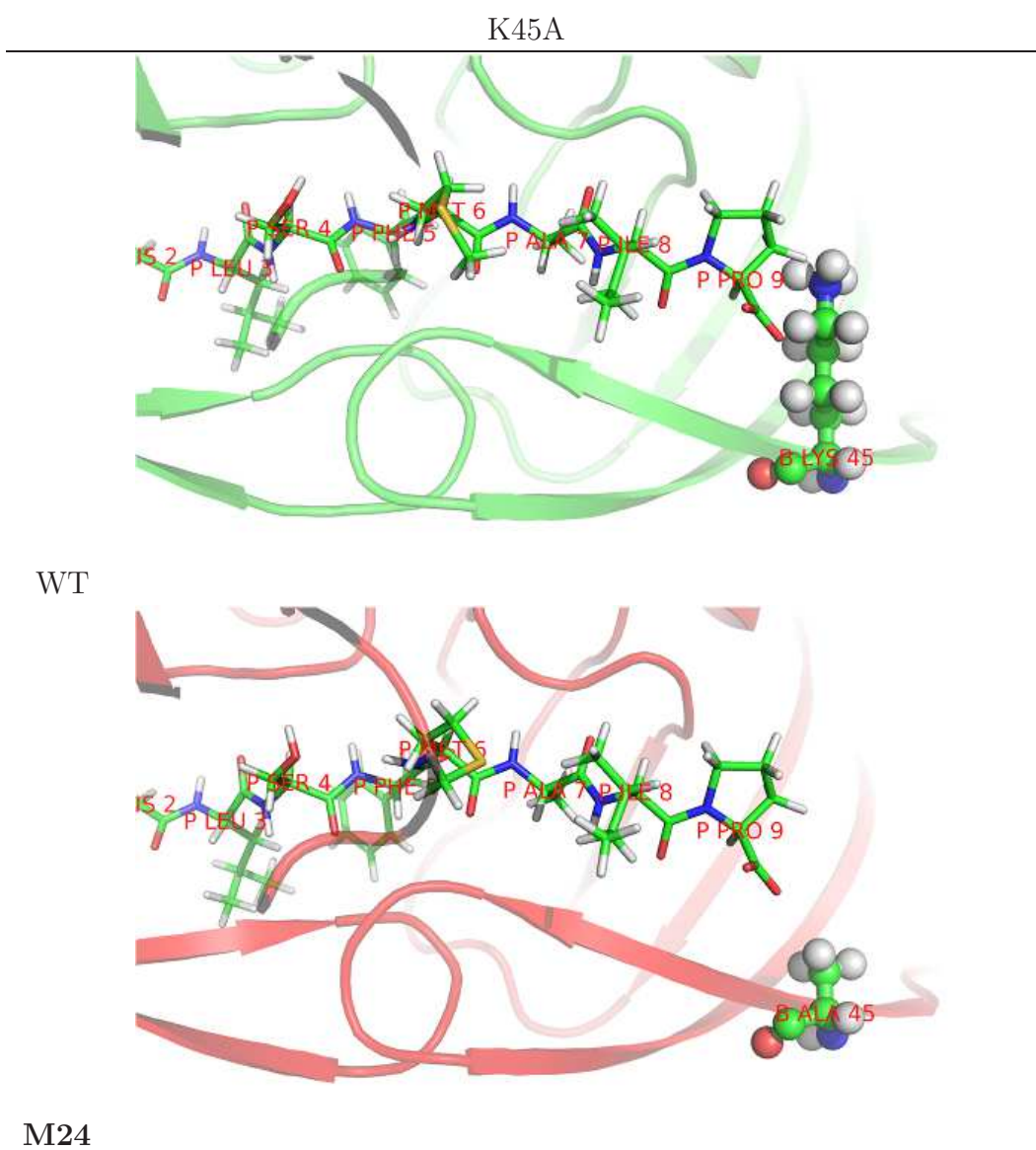


Figure 11S. Chain B, residue 45.



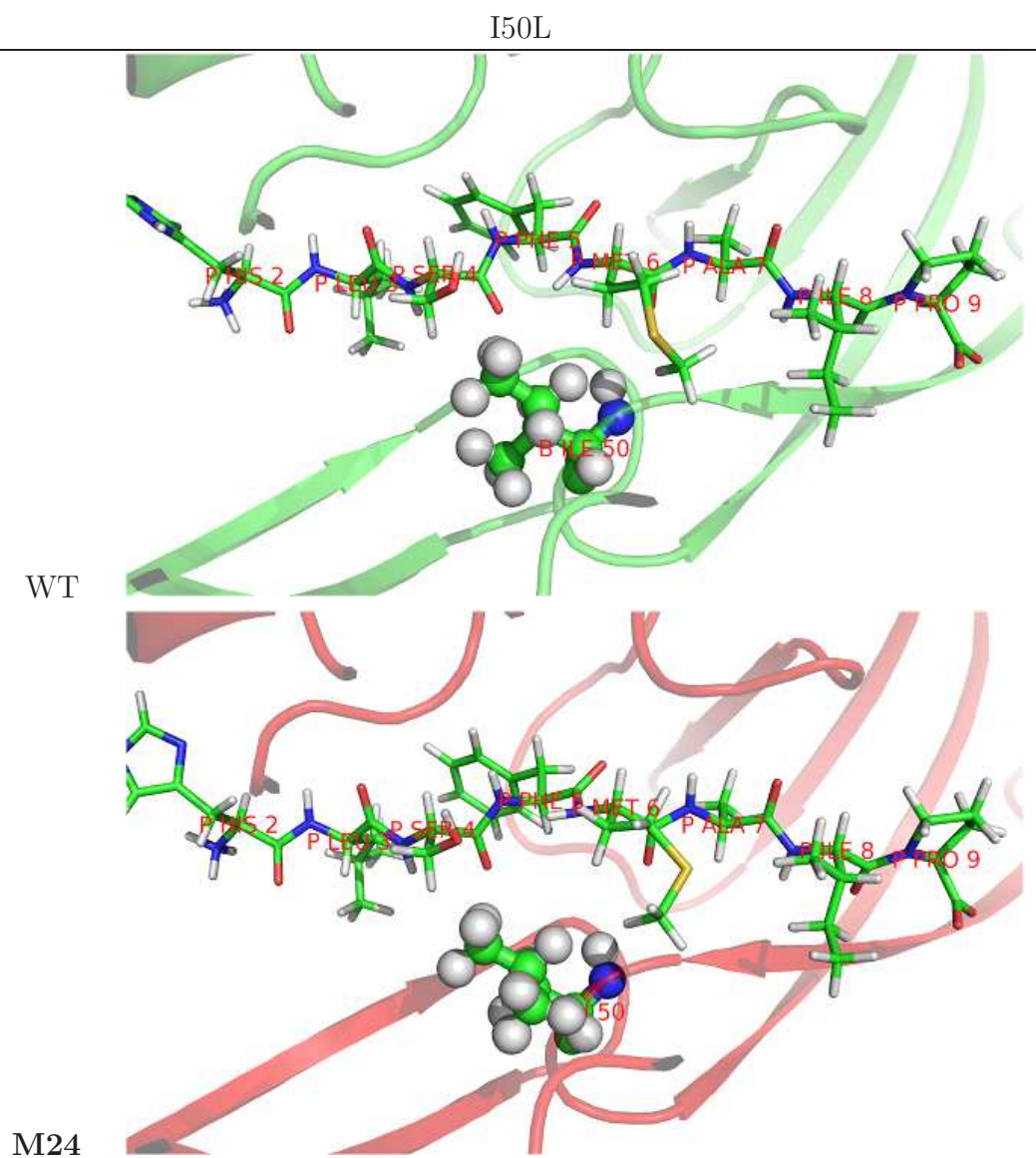


Figure 12S. Chain B, residue 50.

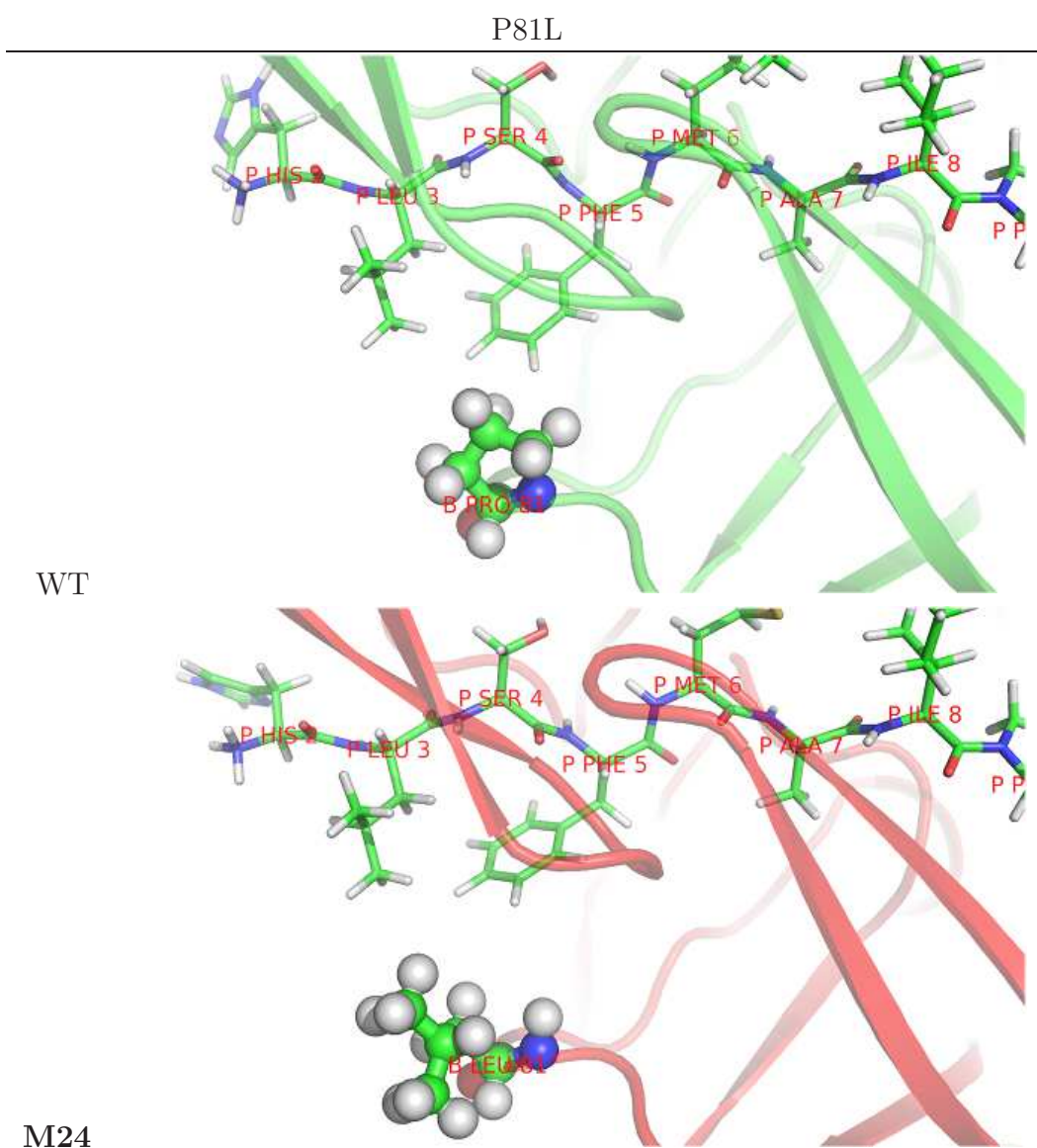


Figure 13S. Chain B, residue 81.

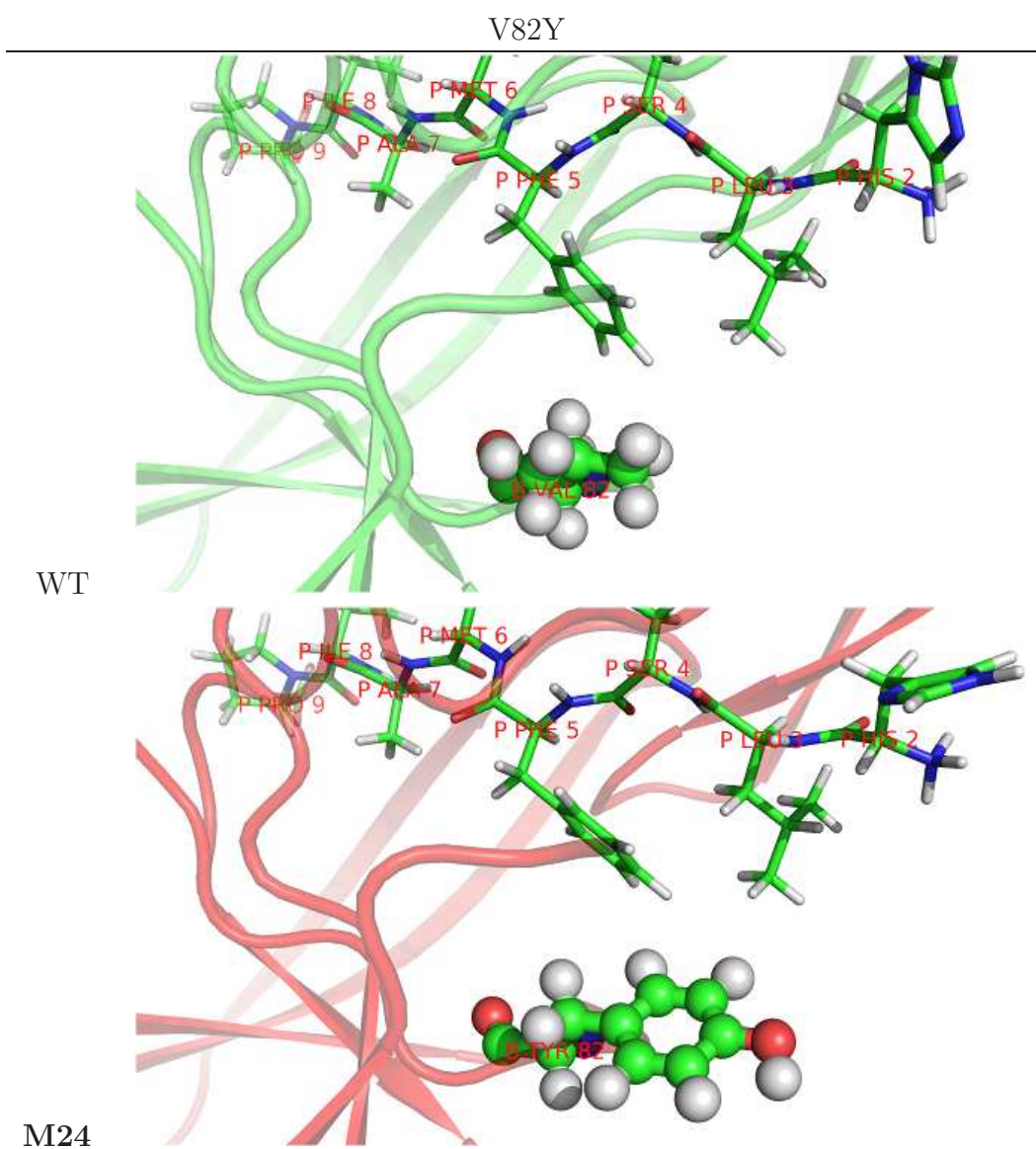


Figure 14S. Chain B, residue 82.

## Title

Author1<sup>1</sup>, Author2<sup>2</sup>, Author3<sup>3,\*</sup>

1 Author1 Dept/Program/Center, Institution Name, City, State, Country

2 Author2 Dept/Program/Center, Institution Name, City, State, Country

3 Author3 Dept/Program/Center, Institution Name, City, State, Country

\* E-mail: Corresponding author@institute.edu

## Abstract

## Author Summary

## Introduction

## Results

### Subsection 1

### Subsection 2

## Discussion

## Materials and Methods

## Acknowledgments

## Figure Legends

## Tables